

A

01/18/00

- ☒ a. ☒ Fees required under 37 CFR 1.16 (National filing fees).
- ☒ b. ☒ Fees required under 37 CFR 1.17 (National application processing fees).
- ☒ A check in the amount of \$ 345 is enclosed.
- ☐ The above filing fee will be paid along with Applicant(s) Response to the Notice to File Missing Parts.

2. ☒ Specification, Total Pages 166
3. ☒ 62 Sheets of Formal Drawing(s) (35 USC 113)
4. ☒ Declaration and Power of Attorney; [Total Pages 2]
- a. ☐ Newly executed (original or copy)
- b. ☒ Copy from a prior application (37 CFR 1.63(d))
(for continuation/divisional with Box 16 completed)
- i. ☐ DELETION OF INVENTOR(S) Signed statement
attached deleting inventor(s) named in the prior
application, see 37 CFR 1.63(d)(2) & 1.33(b).
5. ☐ Microfiche Computer Program (Appendix)
6. ☒ Nucleotide and/or Amino Acid Sequence Submission (if applicable, all necessary)
- a. ☐ Computer Readable Copy
- b. ☒ Paper Copy (identical to computer copy)
- c. ☒ Statement verifying identity of above paper copy with
computer readable copy in prior application

ACCOMPANYING APPLICATION PARTS

7. ☐ Assignment Papers (cover sheet & document(s) (including \$40.00 fee)
8. ☒ 37 CFR 3.73(b) Statement (when there is an assignee); ☒ Power of Attorney
9. ☐ English Translation Document (if applicable)
10. ☒ Information Disclosure Statement (IDS)/PTO-1449
11. ☒ Preliminary Amendment
12. ☒ Return Receipt Postcard (MPEP 503) (Should be specifically itemized)
13. ☒ Small Entity Statement(s)
☒ Statement as filed in prior application; status still proper and desired.
14. ☐ Certified Copy of Priority Document(s) (if foreign priority is claimed)
Foreign Priority is _____
15. ☐ Other: _____
16. **If a CONTINUING APPLICATION**, check appropriate box and supply the requisite
information below and in a preliminary amendment:

- ☐ Continuation ☒ Divisional ☐ Continuation in Part (CIP)
of prior Application No: 09/276,820; Filed March 26, 1999, which is a CIP of

09/263,814 filed 3/8/99, which is a CIP of 09/253,022 filed 2/19/99, which is a CIP of 09/159,643 filed 9/24/98,
which is a CIP of 08,941,223 filed 09/26/97

Prior Application Information: Examiner Ram Shukla

Group/Art Unit: 1632

For CONTINUATION or DIVISIONAL APPS only: The entire disclosure of the prior application, from which an oath or declaration is supplied under Box 4b, is considered a part of the disclosure of the accompanying continuation or divisional application and is hereby incorporated by reference. The incorporation can only be relied upon when a portion has been inadvertently omitted from the submitted application parts.

17. **CORRESPONDENCE ADDRESS**

Customer Number or Bar Code Label **000826**

Attention Of: Anne Brown

Signature: _____

Anne Brown

Attorney of Record: Anne Brown

Attorney Registration No. 36,463

Tel Raleigh Office (919) 420-2200

Fax Raleigh Office (919) 420-2260

ALSTON & BIRD LLP

P.O. Drawer 34009

Charlotte NC 28234-4009

"Express Mail" mailing label number EL247263380US

Date of Deposit January 18, 2000

I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on the date indicated above and is addressed to Box Patent Application, Assistant Commissioner For Patents, Washington, DC 20231.

Nora C Martinez

Nora C. Martinez

RTA01/2071689v1

Statement Claiming Small Entity Status
(37 C.F.R. §§ 1.9(d) and 1.27(c)) -- Small Business Concern

Applicant or Patentee: John J. HARRINGTON, Bruce SHERF and Stephen RUNDLETTAppl. or Patent No.: 09/276,820Attorney Docket No. 1522.0030004/MAC/BJDFiled or Issued: March 26, 1999For: Compositions and Methods for Non-targeted Activation of Endogenous Genes

I hereby state that I am

- ☐ the owner of the small business concern identified below:
☒ an official of the small business concern empowered to act on behalf of the concern identified below:

NAME OF SMALL BUSINESS CONCERN Athersys, Inc.ADDRESS OF SMALL BUSINESS CONCERN 11000 Cedar Avenue, Cleveland, Ohio 44106

I hereby state that the above identified small business concern qualifies as a small business concern as defined in 13 C.F.R. § 121.3-18, and reproduced in 37 C.F.R. § 1.9 (d), for purposes of paying reduced fees under section 41(a) and (b) of Title 35, United States Code, in that the number of employees of the concern, including those of its affiliates, does not exceed 500 persons. For purposes of this statement, (1) the number of employees of the business concern is the average over the previous fiscal year of the concern of the persons employed on a full-time, part-time or temporary basis during each of the pay periods of the fiscal year, and (2) concerns are affiliates of each other when either, directly or indirectly, one concern controls or has the power to control the other, or a third party or parties controls or has the power to control both.

I hereby state that rights under contract or law have been conveyed to and remain with the small business concern identified above with regard to the invention described in:

- ☐ the specification filed herewith with title as listed above.
☒ the application identified above.
☐ the patent identified above.

If the rights held by the above identified small business concern are not exclusive, each individual, concern or organization having rights in the invention must file separate statements indicating their status as small entities, and no rights to the invention are held by any person, other than the inventor, who would not qualify as an independent inventor under 37 C.F.R. § 1.9(c) if that person made the invention or by any concern which would not qualify as a small business concern under 37 C.F.R. § 1.9(d) or a nonprofit organization under 37 C.F.R. § 1.9(e).

Each person, concern or organization having any rights in the invention (other than the small business concern named above) is listed below:

- ☒ no such person, concern, or organization exists.
☐ each person, concern, or organization is listed below.

NAME _____

ADDRESS _____

☐ INDIVIDUAL☒ SMALL BUSINESS CONCERN☐ NONPROFIT ORGANIZATION

Separate statements are required from each named person, concern or organization having rights to the invention averring to their status as small entities. (37 C.F.R. § 1.27)

I acknowledge the duty to file, in this application or patent, notification of any change in status resulting in loss of entitlement to small entity status prior to paying, or at the time of paying, the earliest of the issue fee or any maintenance fee due after the date on which status as a small entity is no longer appropriate. (37 C.F.R. § 1.28(b))

NAME OF PERSON SIGNING _____

TITLE IN ORGANIZATION _____

ADDRESS OF PERSON SIGNING _____

SIGNATURE _____

DATE _____

09/276,820-0110000

Attorney's Docket No. 5817-7L

PATENT

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re:	Harrington, <i>et al.</i>	Group Art Unit:	Not Yet Assigned
Appl. No.:	Not Yet Assigned	Examiner:	Not Yet Assigned
Filed:	Concurrently Herewith		
For:	COMPOSITIONS AND METHODS FOR NON-TARGETED ACTIVATION OF ENDOGENOUS GENES		

January 18, 2000

Assistant Commissioner for Patents
Washington, DC 20231

PRELIMINARY AMENDMENT

Dear Sir:

Please amend the above-identified application as follows:

In The Specification:

Below the heading "Cross-Reference to Related Applications" and after the words "This application", insert the following --is a divisional application of U.S. Application No. 09/276,820, filed March 26, 1999, entitled "COMPOSITIONS AND METHODS FOR NON-TARGETED ACTIVATION OF ENDOGENOUS GENES" which--; and on line 2, in the blank, please insert --09/263,814--.

In The Claims:

Please cancel claims 1-57, without prejudice to or disclaimer of the subject matter contained therein.

00443371-011800

58. (New) A method for drug discovery comprising:

59. (New) A method for drug discovery comprising:

60. (New) The method of claim 59, further comprising concentrating said cell-conditioned media prior to said screening in (c).

61. (New) The method of claim 59, further comprising isolating said gene product prior to said screening in (c).

REMARKS

No new matter has been added by the foregoing amendment to the specification, which has been made solely to provide the application number for the priority application filed on March 26, 1999, which number was not available at the time of filing of the present application.

The foregoing amendments to the claims are fully supported in the specification as originally filed. Specifically, support for new claims 58-61 may be found in the specification at pages 6-17, at pages 38-44, at pages 50-53, at pages 57-61, at pages 68-118, and throughout the Examples. Accordingly, the foregoing amendments to the claims do not add new matter; their entry is therefore respectfully requested. Upon entry of the foregoing amendments, claims 58-61 are pending in the present application.

Applicants believe that the present application is now in condition for examination. If the Examiner believes, for any reason, that personal communication will expedite prosecution of this application, the Examiner is invited to telephone the undersigned at the number provided.

Prompt and favorable consideration of the foregoing amendments, and entry of the same into the present application, are respectfully requested.

It is not believed that extensions of time or fees for net addition of claims are required, beyond those that may otherwise be provided for in documents accompanying this paper. However, in the event that additional extensions of time are necessary to allow consideration of

In re: Harrington *et al.*
Appl. No.: Not Yet Assigned
Filed: Concurrently Herewith
Page 4

this paper, such extensions are hereby petitioned under 37 CFR § 1.136(a), and any fee required therefore (including fees for net addition of claims) is hereby authorized to be charged to Deposit Account No. 16-0605.

Respectfully submitted,



Anne Brown
Registration No. 36,463

ALSTON & BIRD LLP

P.O. Drawer 34009
Charlotte, NC 28234
Tel Raleigh Office (919) 420-2200
Fax Raleigh Office (919) 420-2260

"Express Mail" Mailing Label Number EL247263380US
Date of Deposit: January 18, 2000

I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on the date indicated above and is addressed to Box Patent Application, Assistant Commissioner for Patents, Washington, DC 20231.



Nora C. Martinez

Compositions and Methods for Non-targeted Activation of Endogenous Genes

CROSS REFERENCE TO RELATED APPLICATIONS

5 This application is a continuation-in-part of U.S. Application No. _____ of John J. Harrington, Bruce Sherf, and Stephen Rundlett, entitled "Compositions and Methods for Non-targeted Activation of Endogenous Genes," filed March 8, 1999, which is a continuation-in-part of U.S. Application No. 09/253,022, filed February 19, 1999, which is a continuation-in-part of U.S. Application No. 09/159,643, filed September 24, 1998, which is a continuation-in-part of U.S. Application No. 08/941,223, filed September 26, 1997, the disclosures of all of which are incorporated herein by reference in their entireties.

BACKGROUND OF THE INVENTION

Field of the Invention

15 The present invention is in the fields of molecular biology and cellular biology. The invention is directed generally to activation of gene expression or causing over-expression of a gene by recombination methods *in situ*. More specifically, the invention is directed to activation of endogenous genes by non-targeted integration of specialized activation vectors, which are provided by the invention, into the genome of a host cell. The invention also is directed to
20 methods for the identification, activation, and isolation of genes that were heretofore undiscoverable, and to host cells and vectors comprising such isolated genes. The invention also is directed to isolated genes, gene products, nucleic acid molecules, and compositions comprising such genes, gene products and nucleic acid molecules, that may be used in a variety of therapeutic and diagnostic
25 applications. Thus, by the present invention, endogenous genes, including those associated with human disease and development, may be identified, activated, and

isolated without prior knowledge of the sequence, structure, function, or expression profile of the genes.

Related Art

Identification and over-expression of novel genes associated with human disease is an important step towards developing new therapeutic drugs. Current approaches to creating libraries of cells for protein over-expression are based on the production and cloning of cDNA. Thus, in order to identify a new gene using this approach, the gene must be expressed in the cells that were used to make the library. The gene also must be expressed at sufficient levels to be adequately represented in the library. This is problematic because many genes are expressed only in very low quantities, in a rare population of cells, or during short developmental periods.

Furthermore, because of the large size of some mRNAs, it is difficult or impossible to produce full length cDNA molecules capable of expressing the biologically active protein. Lack of full-length cDNA molecules has also been observed for small mRNAs and is thought to be related to sequences in the message that are difficult to produce by reverse transcription or that are unstable during propagation in bacteria. As a result, even the most complete cDNA libraries express only a fraction of the entire set of possible genes.

Finally, many cDNA libraries are produced in bacterial vectors. Use of these vectors to express biologically active mammalian proteins is severely limited since most mammalian proteins do not fold correctly and/or are improperly glycosylated in bacteria.

Therefore, a method for creating a more representative library for protein expression, capable of facilitating faithful expression of biologically active proteins, would be extremely valuable.

Current methods for over-expressing proteins involve cloning the gene of interest and placing it, in a construct, next to a suitable promoter/enhancer,

polyadenylation signal, and splice site, and introducing the construct into an appropriate host cell.

An alternative approach involves the use of homologous recombination to activate gene expression by targeting a strong promoter or other regulatory sequence to a previously identified gene.

WO 90/14092 describes *in situ* modification of genes, in mammalian cells, encoding proteins of interest. This application describes single-stranded oligonucleotides for site-directed modification of genes encoding proteins of interest. A marker may also be included. However, the methods are limited to providing an oligonucleotide sequence substantially homologous to a target site. Thus, the method requires knowledge of the site required for activation by site-directed modification and homologous recombination. Novel genes are not discoverable by such methods.

WO 91/06667 describes methods for expressing a mammalian gene *in situ*. With this method, an amplifiable gene is introduced next to a target gene by homologous recombination. When the cell is then grown in the appropriate medium, both the amplifiable gene and the target gene are amplified and there is enhanced expression of the target gene. As above, methods of introducing the amplifiable gene are limited to homologous recombination, and are not useful for activating novel genes whose sequence (or existence) is unknown.

WO 91/01140 describes the inactivation of endogenous genes by modification of cells by homologous recombination. By these methods, homologous recombination is used to modify and inactivate genes and to produce cells which can serve as donors in gene therapy.

WO 92/20808 describes methods for modifying genomic target sites *in situ*. The modifications are described as being small, for example, changing single bases in DNA. The method relies upon genomic modification using homologous DNA for targeting.

WO 92/19255 describes a method for enhancing the expression of a target gene, achieved by homologous recombination in which a DNA sequence is

integrated into the genome or large genomic fragment. This modified sequence can then be transferred to a secondary host for expression. An amplifiable gene can be integrated next to the target gene so that the target region can be amplified for enhanced expression. Homologous recombination is necessary to this targeted approach.

WO 93/09222 describes methods of making proteins by activating an endogenous gene encoding a desired product. A regulatory region is targeted by homologous recombination and replacing or disabling the region normally associated with the gene whose expression is desired. This disabling or replacement causes the gene to be expressed at levels higher than normal.

WO 94/12650 describes a method for activating expression of and amplifying an endogenous gene *in situ* in a cell, which gene is not expressed or is not expressed at desired levels in the cell. The cell is transfected with exogenous DNA sequences which repair, alter, delete, or replace a sequence present in the cell or which are regulatory sequences not normally functionally linked to the endogenous gene in the cell. In order to do this, DNA sequences homologous to genomic DNA sequences at a preselected site are used to target the endogenous gene. In addition, amplifiable DNA encoding a selectable marker can be included. By culturing the homologously recombinant cells under conditions that select for amplification, both the endogenous gene and the amplifiable marker are co-amplified and expression of the gene increased.

WO 95/31560 describes DNA constructs for homologous recombination. The constructs include a targeting sequence, a regulatory sequence, an exon, and an unpaired splice donor site. The targeting is achieved by homologous recombination of the construct with genomic sequences in the cell and allows the production of a protein *in vitro* or *in vivo*.

WO 96/29411 describes methods using an exogenous regulatory sequence, an exogenous exon, either coding or non-coding, and a splice donor site introduced into a preselected site in the genome by homologous recombination. In this application, the introduced DNA is positioned so that the transcripts under

control of the exogenous regulatory region include both the exogenous exon and endogenous exons present in either the thrombopoietin, DNase I, or β -interferon genes, resulting in transcripts in which the exogenous and endogenous exons are operably linked. The novel transcription units are produced by homologous recombination.

U.S. Patent No. 5,272,071 describes the transcriptional activation of transcriptionally silent genes in a cell by inserting a DNA regulatory element capable of promoting the expression of a gene normally expressed in that cell. The regulatory element is inserted so that it is operably linked to the normally silent gene. The insertion is accomplished by means of homologous recombination by creating a DNA construct with a segment of the normally silent gene (the target DNA) and the DNA regulatory element used to induce the desired transcription.

U.S. Patent No. 5, 578,461 discusses activating expression of mammalian target genes by homologous recombination. A DNA sequence is integrated into the genome or a large genomic fragment to enhance the expression of the target gene. The modified construct can then be transferred to a secondary host. An amplifiable gene can be integrated adjacent to the target gene so that the target region is amplified for enhanced expression.

Both of the above approaches (construction of an over-expressing construct by cloning or by homologous recombination *in vivo*) require the gene to be cloned and sequenced before it can be over-expressed. Furthermore, using homologous recombination, the genomic sequence and structure must also be known.

Unfortunately, many genes have not yet been identified and/or sequenced. Thus, a method for over-expressing a gene of interest, whether or not it has been previously cloned, and whether or not its sequence and structure are known, would be useful.

BRIEF SUMMARY OF THE INVENTION

The invention is, therefore, generally directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell. The method does not require previous knowledge of the sequence of the endogenous gene or even of the existence of the gene. Hence, the invention is directed to non-targeted gene activation, which as used herein means the activation of endogenous genes by non-targeted or non-homologous (as opposed to targeted or homologous) integration of specialized activation vectors into the genome of a host cell.

The invention also encompasses novel vector constructs for activating gene expression or over-expressing a gene through non-homologous recombination. The novel construct lacks homologous targeting sequences. That is, it does not contain nucleotide sequences that target host cell DNA and promote homologous recombination at the target site, causing over-expressing of a cellular gene via the introduced transcriptional regulatory sequence.

Novel vector constructs include a vector containing a transcriptional regulatory sequence operably linked to an unpaired splice donor sequence and further contains one or more amplifiable markers.

Novel vector constructs include constructs with a transcriptional regulatory sequence operably linked to a translational start codon, a signal secretion sequence, and an unpaired splice donor site; constructs with a transcriptional regulatory sequence, operably linked to a translation start codon, an epitope tag, and an unpaired splice donor site; constructs containing a transcriptional regulatory sequence operably linked to a translational start codon, a signal sequence and an epitope tag, and an unpaired splice donor site; constructs containing a transcriptional regulatory sequence operably linked to a translation

start codon, a signal secretion sequence, an epitope tag, and a sequence-specific protease site, and an unpaired splice donor site.

The vector construct can contain one or more selectable markers for recombinant host cell selection. Alternatively, selection can be effected by phenotypic selection for a trait provided by the activated endogenous gene product.

These vectors, and indeed any of the vectors disclosed herein, and variants of the vectors that will be readily recognized by one of ordinary skill in the art, can be used in any of the methods described herein to form any of the compositions producible by these methods.

The transcriptional regulatory sequence used in the vector constructs of the invention includes, but is not limited to, a promoter. In preferred embodiments, the promoter is a viral promoter. In highly preferred embodiments, the viral promoter is the cytomegalovirus immediate early promoter. In alternative embodiments, the promoter is a cellular, non-viral promoter or inducible promoter.

The transcriptional regulatory sequence used in the vector construct of the invention may also include, but is not limited to, an enhancer. In preferred embodiments, the enhancer is a viral enhancer. In highly preferred embodiments, the viral enhancer is the cytomegalovirus immediate early enhancer. In alternative embodiments, the enhancer is a cellular non-viral enhancer.

In preferred embodiments of the methods described herein, the vector construct be, or may contain, linear RNA or DNA.

The cell containing the vector may be screened for expression of the gene.

The cell over-expressing the gene can be cultured *in vitro* under conditions favoring the production, by the cell, of desired amounts of the gene product (also referred to interchangeably herein as the "expression product") of the endogenous gene that has been activated or whose expression has been increased. The expression product can then be isolated and purified to use, for example, in protein therapy or drug discovery.

Alternatively, the cell expressing the desired gene product can be allowed to express the gene product *in vivo*. In certain such aspects of the invention, the cell containing a vector construct of the invention integrated into its genome may be introduced into a eukaryote (such as a vertebrate, particularly a mammal, more particularly a human) under conditions favoring the overexpression or activation of the gene by the cell *in vivo* in the eukaryote. In related such aspects of the invention, the cell may be isolated and cloned prior to being introduced into the eukaryote.

The invention is also directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence and one or more amplifiable markers into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell.

The cell containing the vector may be screened for over-expression of the gene.

The cell over-expressing the gene is cultured such that amplification of the endogenous gene is obtained. The cell can then be cultured *in vitro* so as to produce desired amounts of the gene product of the amplified endogenous gene that has been activated or whose expression has been increased. The gene product can then be isolated and purified.

Alternatively, following amplification, the cell can be allowed to express the endogenous gene and produce desired amounts of the gene product *in vivo*.

It is to be understood, however, that any vector used in the methods described herein can include one or more amplifiable markers. Thereby, amplification of both the vector and the DNA of interest (i.e., containing the over-expressed gene) occurs in the cell, and further enhanced expression of the endogenous gene is obtained. Accordingly, methods can include a step in which the endogenous gene is amplified.

5 The invention is also directed to methods for over-expressing an endogenous gene in a cell comprising introducing a vector containing a transcriptional regulatory sequence and an unpaired splice donor sequence into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell.

The cell containing the vector may be screened for expression of the gene.

10 The cell over-expressing the gene can be cultured *in vitro* so as to produce desirable amounts of the gene product of the endogenous gene whose expression has been activated or increased. The gene product can then be isolated and purified.

Alternatively, the cell can be allowed to express the desired gene product *in vivo*.

15 The vector construct can consist essentially of the transcriptional regulatory sequence.

The vector construct can consist essentially of the transcriptional regulatory sequence and one or more amplifiable markers.

The vector construct can consist essentially of the transcriptional regulatory sequence and the splice donor sequence.

20 Any of the vector constructs of the invention can also include a secretion signal sequence. The secretion signal sequence is arranged in the construct so that it will be operably linked to the activated endogenous protein. Thereby, secretion of the protein of interest occurs in the cell, and purification of that protein is facilitated. Accordingly, methods can include a step in which the protein expression product is secreted from the cell.

25 The invention also encompasses cells made by any of the above methods. The invention encompasses cells containing the vector constructs, cells in which the vector constructs have integrated into the cellular genome, and cells which are over-expressing desired gene products from an endogenous gene, over-expression being driven by the introduced transcriptional regulatory sequence.

30

The cells can be isolated and cloned.

The methods can be carried out in any cell of eukaryotic origin, such as fungal, plant or animal. In preferred embodiments, the methods of the invention may be carried out in vertebrate cells, and particularly mammalian cells including but not limited to rat, mouse, bovine, porcine, sheep, goat and human cells, and more particularly in human cells.

A single cell made by the methods described above can over-express a single gene or more than one gene. More than one gene in a cell can be activated by the integration of a single type of construct into multiple locations in the genome. Similarly, more than one gene in a cell can be activated by the integration of multiple constructs (i.e., more than one type of construct) into multiple locations in the genome. Therefore, a cell can contain only one type of vector construct or different types of constructs, each capable of activating an endogenous gene.

The invention is also directed to methods for making the cells described above by one or more of the following: introducing one or more of the vector constructs of the invention into a cell; allowing the introduced construct(s) to integrate into the genome of the cell by non-homologous recombination; allowing over-expression of one or more endogenous genes in the cell; and isolating and cloning the cell. The invention is also directed to cells produced by such methods, which may be isolated cells.

The invention also encompasses methods for using the cells described above to over-express a gene, such as an endogenous cellular gene, that has been characterized (for example, sequenced), uncharacterized (for example, a gene whose function is known but which has not been cloned or sequenced), or a gene whose existence was, prior to over-expression, unknown. The cells can be used to produce desired amounts of an expression product *in vitro* or *in vivo*. If desired, this expression product can then be isolated and purified, for example by cell lysis or by isolation from the growth medium (as when the vector contains a secretion signal sequence).

The invention also encompasses libraries of cells made by the above described methods. A library can encompass all of the clones from a single transfection experiment or a subset of clones from a single transfection experiment. The subset can over-express the same gene or more than one gene, for example, a class of genes. The transfection can have been done with a single construct or with more than one construct.

A library can also be formed by combining all of the recombinant cells from two or more transfection experiments, by combining one or more subsets of cells from a single transfection experiment or by combining subsets of cells from separate transfection experiments. The resulting library can express the same gene, or more than one gene, for example, a class of genes. Again, in each of these individual transfections, a unique construct or more than one construct can be used.

Libraries can be formed from the same cell type or different cell types.

The invention is also directed to methods for making libraries by selecting various subsets of cells from the same or different transfection experiments.

The invention is also directed to methods of using the above-described cells or libraries of cells to over-express or activate endogenous genes, or to obtain the gene expression products of such over-expressed or activated genes. According to this aspect of the invention, the cell or library may be screened for the expression of the gene and cells that express the desired gene product may be selected. The cell can then be used to isolate or purify the gene product for subsequent use. Expression in the cell can occur by culturing the cell *in vitro*, under conditions favoring the production of the expression product of the endogenous gene by the cell, or by allowing the cell to express the gene *in vivo*.

In preferred embodiments of the invention, the methods include a process wherein the expression product is isolated or purified. In highly preferred embodiments, the cells expressing the endogenous gene product are cultured under conditions favoring production of sufficient amounts of gene product for

commercial application, and especially for diagnostic, therapeutic and drug discovery uses.

Any of the methods can further comprise introducing double-strand breaks into the genomic DNA in the cell prior to or simultaneously with vector integration.

The invention also is directed to vector constructs that are useful for activating expression of endogenous genes and for isolating the mRNA and cDNA corresponding to the activated genes.

In one such embodiment, the vector construct may comprise (a) a first transcriptional regulatory sequence operably linked to a first unpaired splice donor sequence; (b) a second transcriptional regulatory sequence operably linked to a second unpaired splice donor sequence; and (c) a linearization site, which may be located between the first and second transcriptional regulatory sequences. According to the invention, when the vector construct is transformed into a host cell and then integrates into the genome of the host cell, the first transcriptional regulatory sequence is preferably in an inverted orientation relative to the orientation of the second transcriptional regulatory sequence. In certain preferred such embodiments, the vector may be rendered linear by cleavage at the linearization site.

In another embodiment, the invention provides a linear vector construct having a 3' end and a 5' end, comprising a transcriptional regulatory sequence operably linked to an unpaired spliced donor site, wherein the transcriptional regulatory sequence is oriented in the linear vector construct in an orientation that directs transcription towards the 3' end or the 5' end of the linear vector construct.

In another embodiment, the invention provides a vector construct comprising, in sequential order, (a) a transcriptional regulatory sequence, (b) an unpaired splice donor site, (c) a rare cutting restriction site, and (d) a linearization site.

In another embodiment, the invention provides a vector construct comprising (a) a first transcriptional regulatory sequence operably linked to a

selectable marker lacking a polyadenylation signal; and (b) a second transcriptional regulatory sequence operably linked to an exon-splice donor site complex, wherein the first transcriptional regulatory sequence is in the same orientation in the vector construct as is the second transcriptional regulatory sequence, and wherein the first transcriptional regulatory sequence is upstream of the second transcriptional regulatory sequence in the vector construct.

In additional embodiments, the invention provides vector constructs comprising a transcriptional regulatory sequence operably linked to a selectable marker lacking a polyadenylation signal, and further comprising an unpaired splice donor site.

In another embodiment, the invention provides vector constructs comprising a first transcriptional regulatory sequence operably linked to a selectable marker lacking a polyadenylation signal, and further comprising a second transcriptional regulatory sequence operably linked to an unpaired splice donor site.

According to the invention, the transcriptional regulatory sequence (or first or second transcriptional regulatory sequence, in vector constructs having more than one transcriptional regulatory sequence) may be a promoter, an enhancer, or a repressor, and is preferably a promoter, including an animal cell promoter, a plant cell promoter, or a fungal cell promoter, most preferably a promoter selected from the group consisting of a CMV immediate early gene promoter, an SV40 T antigen promoter, and a β -actin promoter. Other promoters of animal, plant, or fungal cell origin that may be used in accordance with the invention are known in the art and will be familiar to one of ordinary skill in view of the teachings herein.

The selectable marker used in the vector constructs of the invention may be any marker or marker gene that, upon integration of a vector containing the selectable marker into the host cell genome, permits the selection of a cell containing or expressing the marker gene. Suitable such selectable markers include, but are not limited to, a neomycin gene, a hypoxanthine phosphoribosyl transferase gene, a puromycin gene, a dihydroorotase gene, a glutamine synthetase

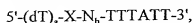
gene, a histidine D gene, a carbamyl phosphate synthase gene, a dihydrofolate reductase gene, a multidrug resistance 1 gene, an aspartate transcarbamylase gene, a xanthine-guanine phosphoribosyl transferase gene, an adenosine deaminase gene, and a thymidine kinase gene.

5 In related embodiments, the invention provides vector constructs comprising a positive selectable marker, a negative selectable marker, and an unpaired splice donor site, wherein the positive and negative selectable markers and the splice donor site are oriented in the vector construct in an orientation that results in expression of the positive selectable marker in active form, and either
10 non-expression of said negative selectable marker or expression of the negative selectable marker in inactive form, when the vector construct is integrated into the genome of a eukaryotic host cell and activates an endogenous gene in the genome. In certain preferred such embodiments, either the positive selection marker, the negative selection marker, or both, may lack a polyadenylation signal. The
15 positive selection marker used in these aspects of the invention may be any selection marker that, upon expression, produces a protein capable of facilitating the isolation of cells expressing the marker, including but not limited to a neomycin gene, a hypoxanthine phosphoribosyl transferase gene, a puromycin gene, a dihydroorotase gene, a glutamine synthetase gene, a histidine D gene, a carbamyl
20 phosphate synthase gene, a dihydrofolate reductase gene, a multidrug resistance 1 gene, an aspartate transcarbamylase gene, a xanthine-guanine phosphoribosyl transferase gene, or an adenosine deaminase gene. Analogously, the negative selection marker used in these aspects of the invention may be any selection marker that, upon expression, produces a protein capable of facilitating removal
25 of cells expressing the marker, including but not limited to a hypoxanthine phosphoribosyl transferase gene, a thymidine kinase gene, or a diphtheria toxin gene.

The invention also is directed to eukaryotic host cells, which may be isolated host cells, comprising one or more of the vector constructs of the invention. Preferred such eukaryotic host cells include, but are not limited to,
30

animal cells (including, but not limited to, mammalian (particularly human) cells, insect cells, avian cells, annelid cells, amphibian cells, reptilian cells, and fish cells), plant cells, and fungal (particularly yeast) cells. In certain such host cells, the vector construct may be integrated into the genome of the host cell.

The invention also is directed to primer molecules comprising a PCR-amplifiable sequence and a degenerate 3' terminus. Primer molecules according to this aspect of the invention preferably have the general structure:



wherein a is a whole number from 1 to 100 (preferably from 10 to 30), X is a PCR-amplifiable sequence consisting of a nucleic acid sequence of about 10-20 nucleotides in length, N is any nucleotide, and b is a whole number from 0 to 6. One preferred such primer has the nucleotide sequence 5'-TTTTTTT-TTTTCGT CAGCGCCGCATC NNNNTTATT-3' (SEQ ID NO:10). In related embodiments, the primer molecules according to this aspect of the invention may be biotinylated.

The invention also is directed to methods for first strand cDNA synthesis comprising (a) annealing a first primer of the invention (such as the primer described above) to an RNA template molecule to form a first primer-RNA complex, and (b) treating this first primer-RNA complex with reverse transcriptase and one or more deoxynucleoside triphosphate molecules under conditions favoring the reverse transcription of the first primer-RNA complex to synthesize a first strand cDNA.

The invention also is directed to methods for isolating activated genes, particularly from a host cell genome. These methods of the invention exploit the structure of the mRNA molecules produced using the non-targeted gene activation vectors of the invention. One such method of the invention comprises, for example, (a) introducing a vector construct comprising a transcriptional regulatory sequence and an unpaired splice donor site into a host cell (preferably one of the eukaryotic host cells described above), (b) allowing the vector construct to integrate into the genome of the host cell by non-homologous recombination,

under conditions such that the vector activates an endogenous gene comprising an exon in the genome, (c) isolating RNA from the host cell, (d) synthesizing first strand cDNA according to the method of the invention described above, (e) annealing a second primer specific for the vector-encoded exon to the first strand cDNA to create a second primer-first strand cDNA complex, and (f) contacting the second primer-first strand cDNA complex with a DNA polymerase under conditions favoring the production of a second strand cDNA substantially complementary to the first strand cDNA. Methods according to this aspect of the invention may comprise one or more additional steps, such as treating the second strand cDNA with a restriction enzyme that cleaves at a restriction site located on the vector downstream of the unpaired splice donor site, or amplifying the second strand cDNA using a third primer specific for the vector-encoded exon and a fourth primer specific for the second primer. The invention also is directed to isolated genes produced according to these methods, and to vectors (which may be expression vectors) and host cells comprising these isolated genes. The invention also is directed to methods of producing a polypeptide, comprising cultivating a host cell comprising the isolated gene (or a vector, particularly an expression vector, comprising the isolated gene), and culturing the host cell under conditions favoring the expression by the host cell of a polypeptide encoded by the isolated gene. The invention also provides additional methods of producing a polypeptide, comprising introducing into a host cell a vector comprising a transcriptional regulatory sequence operably linked to an exonic region followed by an unpaired splice donor site, and culturing the host cell under conditions favoring the expression by said host cell of a polypeptide encoded by the exonic region, wherein the exon contains a translational start site positioned at any of the open reading frame positions relative to the 5'-most base of the unpaired splice donor site (*e.g.*, the "A" in the ATG start codon may be at position -3 or at an increment of 3 bases upstream therefrom (*e.g.*, -6, -9, -12, -15, -18, etc.), at position -2 or at an increment of 3 bases upstream therefrom (*e.g.*, -5, -8, -11, -14, -17, -20, etc.), or at position -1 or at an increment of 3 bases upstream

therefrom (e.g., -4, -7, -10, -13, -16, -19, etc.), relative to the 5'-most base of the splice donor site). In related embodiments, the methods of the invention may further comprise isolating the polypeptide. The invention also is directed to polypeptides, which may or may not be isolated polypeptides, produced according to these methods.

Other preferred embodiments of the present invention will be apparent to one of ordinary skill in light of the following drawings and description of the invention, and of the claims.

BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1. Schematic diagram of gene activation events described herein.

The activation construct is transfected into cells and allowed to integrate into the host cell chromosomes at DNA breaks. If breakage occurs upstream of a gene of interest (e.g., Epo), and the appropriate activation construct integrates at the break such that its regulatory sequence becomes operably linked to the gene of interest, activation of the gene will occur. Transcription and splicing produce a chimeric RNA molecule containing exonic sequences from the activation construct and from the endogenous gene. Subsequent translation will result in the production of the protein of interest. Following isolation of the recombinant cell, gene expression can be further enhanced via gene amplification.

FIG. 2. Schematic diagram of non-translated activation constructs. The arrows denote promoter sequences. The exonic sequences are shown as open boxes and the splice donor sequence is indicated by S/D. Construct numbers corresponding to the description below are shown on the left. The selectable and amplifiable markers are not shown.

FIG. 3. Schematic diagram of translated activation constructs. The arrows denote promoter sequences. The exonic sequences are shown as open boxes and the splice donor sequence is indicated by S/D. The translated, signal peptide, epitope tag, and protease cleavage sequences are shown in the legend

below the constructs. Construct numbers corresponding to the description below are shown on the left. The selectable and amplifiable markers are not shown.

FIG. 4. Schematic diagram of an activation construct capable of activating endogenous genes.

FIG. 5A-5D. Nucleotide sequence of pRIG8R1-CD2 (SEQ ID NO:7).

FIG. 6A-6C. Nucleotide sequence of pRIG8R2-CD2 (SEQ ID NO:8).

FIG. 7A-7C. Nucleotide sequence of pRIG8R3-CD2 (SEQ ID NO:9).

FIG. 8A-8F. Examples of poly(A) trap vectors. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Untranslated regions are designated by hatched boxes and open reading frames are designated by open boxes. The following designations were used: splice donor site (S/D), signal secretion sequence (SP), epitope tag (ET), neomycin resistance gene (Neo). In the vectors depicted in Fig. 8B-8E, it is possible to omit the splice donor site immediately downstream of the Neo gene. In vectors lacking a splice donor site between the neo gene and the downstream promoter, the Neo transcript will utilize the splice donor site located 3' of the downstream promoter. In addition, as shown in the vectors depicted in Fig. 8B-8E, a downstream promoter may drive expression of an exon. It is recognized that this exon, when present, may encode codons in any reading frame. Using multiple vectors, codons in each of the 3 possible reading frames can be created.

FIG. 9A-9F. Examples of splice acceptor trap vectors containing a positive and a negative selectable marker driven from a single promoter. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Untranslated regions are designated by hatched boxes. Poly(A) signals are not present in these examples.

As described in the specification, however, poly(A) signals may be placed on the vector 3' of either or both selectable markers. The following designations were used: splice donor site (S/D), signal secretion sequence (SP), epitope tag (ET), internal ribosome entry site (ires), hypoxanthine phosphoribosyl transferase (HPRT), and neomycin resistance gene (Neo). In these examples, Neo represents the positive selectable marker and HPRT represents the negative selectable marker. In the vectors shown in Fig. 9C and 9F, the region designated exon contains a translation start codon. As described in the Detailed Description, the exon may encode a methionine residue, a partial signal sequence, a full signal secretion sequence, a portion of a protein, or an epitope tag. In addition, the codons may be present in any reading frame relative to the splice donor site. In other vector examples not shown, the region designated exon lacks a translation start codon.

FIG. 10A-10F. Examples of splice acceptor trap vectors containing a positive and negative selectable marker driven from different promoters. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Untranslated regions are designated by hatched boxes. Poly(A) signals are not present in these examples. As described in the specification, however, poly(A) signals may be placed on the vector 3' of either or both selectable markers. The following designations were used: splice donor site (S/D), internal ribosome entry site (ires), hypoxanthine phosphoribosyl transferase (HPRT), and neomycin resistance gene (Neo). In the vectors shown in Figs. 10A-10F, Neo represents the positive selectable marker and HPRT represents the negative selectable marker. As shown, the vectors depicted in Figs. 10A-10F do not contain a splice donor site 3' of the Neo gene; however, in other vectors not shown, a splice donor site may be located 3' of the Neo gene to facilitate splicing of the positive selection marker to an endogenous exon. In the vectors shown in Fig. 10C and 10F, the region designated exon

contains a translation start codon. As described in the Detailed Description, the exon may encode a methionine residue, a partial signal sequence, a full signal secretion sequence, a portion of a protein, or an epitope tag. In addition, the codons may be present in any reading frame relative to the splice donor site. In other vector examples not shown, the region designated exon lacks a translation start codon.

FIG. 11A-11C. Schematic diagram of bidirectional activation vectors.

The arrows denote promoter sequences. The exons are shown as checkered boxes and splice donor sites are indicated by S/D. The hatched boxes indicate exon sequences operably linked to the upstream promoter. It is understood that the exons on these vectors may be untranslated, or may contain a start codon and additional codons as described herein. As illustrated in the vectors depicted in Fig. 11B-11C, the vectors may contain a selectable marker. In these vectors, the neomycin resistance (Neo) gene is illustrated. In Fig. 11B, a polyadenylation signal (pA) is located downstream of the selectable marker. In Fig. 11C, polyadenylation signals are omitted from the vector.

FIG. 12A-12G. Examples of vectors useful for recovering exon I from

activated endogenous genes. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Untranslated regions are designated by hatched boxes. Poly(A) signals are not present in the vectors depicted. As discussed in the Detailed Description, however, poly(A) signals may be placed on the vector 3' of either or both selectable markers. The following designations were used: splice donor site (S/D), internal ribosome entry site (ires), hypoxanthine phosphoribosyl transferase (HPRT), and neomycin resistance gene (Neo). In these examples, Neo represents the positive selectable marker and HPRT represents the negative selectable marker. It is also recognized that in these examples, the region designated exon, when present, lacks a translation start codon. In other examples

not shown, the region designated exon contains a translation start codon. Furthermore, when the vector exon contains a translation start codon, the exon may encode a methionine residue, a partial signal sequence, a full signal secretion sequence, a portion of a protein, or an epitope tag. In addition, the codons may be present in each reading frame relative to the splice donor site.

FIG. 13. Illustration depicting two transcripts produced from the integrated vectors described in Figures 12A-12G. DNA strands are depicted as horizontal lines. Vector DNA is shown as a black line. Endogenous genomic DNA is shown as a grey line. Rectangles depict exons. Vector-encoded exons are shown as open rectangles, while endogenous exons are shown as shaded boxes. S/D denotes a splice donor site. Following integration, the vector encoded promoters activate transcription of the endogenous gene. Transcription resulting from the upstream promoter produces a spliced RNA molecule containing the vector encoded exon joined to the second and subsequent exons from an endogenous gene. Transcription from the downstream promoter, on the other hand, produces a transcript containing the sequences downstream of the integrated joined to exon I and the subsequent exons from an endogenous gene.

FIG. 14A-14B. Nucleotide sequence of pRIG1 (SEQ ID NO:18).

FIG. 15A-15B. Nucleotide sequence of pRIG21b (SEQ ID NO:19).

FIG. 16A-16B. Nucleotide sequence of pRIG22b (SEQ ID NO:20).

FIG. 17A-17G. Examples of poly(A) trap vectors. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions. The following designations were used: splice donor site (S/D), signal secretion sequence (SP), epitope tag (ET), neomycin resistance gene (Neo), vector promoter #1 (VP#1), and vector promoter #2 (VP#2). As shown in the vectors depicted in Fig. 17C-17G, a promoter operably linked to an exon and an unpaired splice donor site can be positioned upstream of

the selectable marker. It is recognized that this exon, when present, may encode codons a start codon in any reading frame relative to the splice donor site. To activate protein expression from genes with different reading frames, three separate vectors can be used, each with a start codon in a different reading frame relative to the splice donor site.

FIG. 18. Illustration of the transcripts produced by the vector from Fig. 17C upon integration into a host cell genome upstream of a multi-exon endogenous gene. Each horizontal line represents a DNA molecule. Vertical lines running through the DNA strand mark the upstream and downstream vector/cellular genome boundaries. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions. The endogenous exons are numbered using roman numerals. The following designations were used: splice donor site (S/D), neomycin resistance gene (Neo), vector promoter #1 (VP#1), vector promoter #2 (VP#2), endogenous promoter (EP) and polyadenylation signal (pA). Following integration, vector promoter #1 expresses a chimeric transcript containing the Neo gene linked to the genomic sequences downstream of the integration site, including the processed (spliced) exons from the endogenous gene. Since transcript #1 contains a poly (A) signal from the endogenous gene, the Neo gene product will be efficiently produced, thereby conferring drug resistance on the cell. In addition to transcript #1, the integrated vector will generate a second transcript, designated transcript #2, originating from vector promoter#2. The structure of transcript #2 facilitates efficient translation of the protein encoded by the endogenous gene. As exemplified in Figure 17, vectors containing alternative coding information in the vector encoded exon can be used to produce different chimeric proteins, containing, for example, signal sequences and/or epitope tags.

FIG. 19. Example of dual positive selectable marker vector. The vector is illustrated schematically in its linearized form. The horizontal line represents a

DNA molecule. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions. Poly(A) signals are not present in these examples. The following designations were used: splice donor site (S/D), hygromycin resistance gene (Hyg), neomycin resistance gene (Neo), vector promoter #1, and vector promoter #2.

FIG. 20A-20B. Examples of transcripts produced by a dual positive selectable marker vector integrated into a host cell genome adjacent to an endogenous gene. Figure 20A illustrates the transcripts produced upon vector integration near a multi-exon gene. Figure 20B illustrates the transcripts produced upon vector integration near a single exon gene. Each horizontal line represents a DNA molecule. Vertical lines running through the DNA strand mark the upstream and downstream vector/cellular genome boundaries. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of each promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions. The endogenous exons are numbered using roman numerals. The following designations were used: splice donor site (S/D), hygromycin resistance gene (Hyg), neomycin resistance gene (Neo), vector promoter #1 (VP#1), vector promoter #2 (VP#2), endogenous promoter (EP), and polyadenylation signal (pA). Following integration, vector promoter #1 expresses a chimeric transcript containing the Hyg gene linked to the genomic sequences downstream of the integration site, including the processed (spliced) exons from the endogenous gene. Since transcript #1 contains a poly (A) signal from the endogenous gene, the Hyg gene product will be efficiently produced, thereby conferring drug resistance on the cell. In addition to transcript #1, the integrated vector will generate a second transcript, designated transcript #2, originating from vector promoter #2. In figure 20A, the neo gene is removed from transcript #2 upon splicing from the vector encoded splice donor site, and the first endogenous splice

004484331-011800

acceptor located downstream of the vector integration site (i.e. exon II in this example). Since multi-exon genes contain splice acceptor sites at the 5' end of each exon (except exon I), the neo gene will be removed from transcript #2 in cells in which the vector has integrated near, and transcriptionally activated, a multi-exon gene. As a result, cells having activated multi-exon genes may be eliminated by selecting with G418 and hygromycin. In figure 20B, the neo gene is not removed from transcript #2 by splicing, since single exon genes do not contain any splice acceptor sequences. Thus, cells containing a vector integrated near single exon genes will survive double selection with G418 and hygromycin. These cells can be used to efficiently isolate the activated single exon genes using methods described herein.

FIG. 21A-21B. Examples of dual trap vectors containing a positive and a negative selectable marker. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of a promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions. The following designations were used: splice donor site (S/D), hypoxanthine phosphoribosyl transferase (HPRT), neomycin resistance gene (Neo), vector promoter #1 (VP #1), vector promoter #2 (VP#2), and vector promoter #3 (VP#3). In the vectors shown in Figs. 21A-21B, Neo represents the positive selectable marker and HPRT represents the negative selectable marker. In re 21B a third promoter is located upstream of the selectable markers. This upstream promoter is operably linked to an exon and unpaired splice donor site. Fig. The region designated exon contains a translation start codon in this example. As described herein, the exon may encode a methionine residue, a partial signal sequence, a full signal secretion sequence, a portion of a protein, or an epitope tag. In addition, the codons may be present in any reading frame relative to the splice donor site. In other vector examples not shown, the region designated exon lacks a translation start codon.

FIG. 22. Examples of transcripts produced by a dual positive/negative selectable marker vector integrated into a host cell genome upstream of a multi-exon endogenous gene. Each horizontal line represents a DNA molecule. Vertical lines running through the DNA strand mark the upstream and downstream vector/cellular genome boundaries. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of each promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions. The endogenous exons are numbered using roman numerals. The following designations were used: splice donor site (S/D), neomycin resistance gene (Neo), vector promoter #1 (VP#1), vector promoter #2 (VP#2), vector promoter #3 (VP#3), polyadenylation signal (pA), and endogenous promoter (EP). Following integration, vector promoter #1 expresses a chimeric transcript containing the Neo gene linked to the genomic sequences downstream of the integration site, including the processed (spliced) exons from the endogenous gene. Since transcript #1 contains a poly (A) signal from the endogenous gene, the Neo gene product will be efficiently produced, thereby conferring drug resistance on the cell. In addition to transcript #1, the integrated vector will generate a second transcript, designated transcript #2, originating from vector promoter #2. In this example, the vector has integrated upstream of a multi-exon gene. Since multi exon genes contain splice acceptor sites at the 5' end of each exon, the HPRT gene will be removed from transcript #2 in cells in which the vector has integrated near, and transcriptionally activated, a multi-exon gene. As a result, cells containing activated multi-exon genes may be isolated by selecting with G418 and 8-Azaguanine 6-Thioguanine (AgThg). Thus, cells containing a vector integrated near single exon genes will survive double selection with G418 and AgThg. These cells can be used to efficiently isolate the activated multi-exon genes using methods described herein. In addition to transcripts #1 and #2, a third transcript, designated transcript #3 is produced from the integrated vector. Transcript #3, originating from vector promoter #3, contains an exonic sequence suitable for

directing protein expression from the endogenous gene. This occurs following splicing from the first splice donor site downstream of promoter #3 to the first downstream splice acceptor site from the endogenous gene. In addition to directing protein expression, transcript #3, and/or transcripts #1 and/or #2, can be isolated for gene discovery purposes using the methods described herein.

FIG. 23A-23D. Example of a multi-Promoter/Activation Exon Vector.

Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences. Boxes indicate exons. Hatched boxes indicate untranslated regions. It is understood that the exons on these vectors may be untranslated, or may contain a start codon and additional codons as described herein. The following designations were used: splice donor site (S/D), vector promoter #1 (VP #1), vector promoter #2 (VP#2), vector promoter #3 (VP #3), and vector promoter #4 (VP#4). Individual vector activation exons are designated A, B, C, and D. Each activation exon may contain a different structure. The structure of each activation exon and its flanking intron are shown below. It is understood, however, that any activation exon described herein, may be used on these vectors, in any combination and/or order, including exons that encode signal sequences, partial signal sequences, epitope tags, proteins, portions of proteins, and protein motifs. Any of the exons may lack a start codon. In addition, while not illustrated in these examples, these vectors may contain a selectable marker and/or an amplifiable marker. The selectable marker may contain a poly (A) signal or a splice donor site. When present, the splice donor site may be located upstream or downstream of the selectable marker. Alternatively, the selectable marker may not be operably linked to a poly (A) signal and/or a splice donor site.

FIG. 24. Examples of transcripts produced from a multi-Promoter/Activation Exon Vector upon integration into a host cell genome upstream of an endogenous gene. Each horizontal line represents a DNA molecule. Vertical lines running through the DNA strand mark the upstream and downstream vector/cellular genome boundaries. The arrows denote promoter

sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of each promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions.. The endogenous exons are numbered using roman numerals. The following designations were used: splice donor site (S/D), vector promoter #1 (VP #1), vector promoter #2 (VP#2), vector promoter #3 (VP #3), vector promoter #4 (VP#4), endogenous promoter (EP), and polyadenylation signal (pA). Individual vector activation exons are designated A, B, C, and D.. Following integration, each vector encoded promoter is capable of producing a different transcript. Each transcript contains a different activation exon joined to the first downstream splice acceptor site from an endogenous gene (exon II in this example). Individual activation exons are designated by (A), (B), (C), or (D). Endogenous exons are designated by (I), (II), (III), or (IV). Generally, the coding sequence and/or reading frames, if present, are different among the activation exons. While four activation exons are illustrated in this example, any number of activation exons may be present on the integrated vector.

FIG. 25A-25D. Examples of activation vectors useful for detection of protein-protein interactions. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences. Boxes indicate exons. Hatched boxes indicate untranslated regions. The following designations were used: splice donor site (S/D), neomycin resistance gene (Neo). It is also recognized that the DNA binding domain and the Activation domain may be encoded in any reading frame (relative to the splice donor site), allowing activation of endogenous genes with different reading frames.

FIG. 26. Schematic illustration depicting one approach to detecting protein-protein interactions using the vectors shown in Figure 25. Each horizontal line represents a DNA molecule. Vertical lines running through the DNA strand mark the upstream and downstream vector/cellular genome boundaries. The arrows denote promoter sequences located on the DNA molecule, and face in the

direction of transcription. Transcribed regions include all sequences located downstream of each promoter. Boxes indicate exons. Hatched boxes indicate untranslated regions. The endogenous exons are numbered using roman numerals. The following designations were used: splice donor site (S/D), binding domain (BD), activation domain (AD), recognition sequence (RS), and polyadenylation signal (pA). The binding domain vector is shown integrated into the genome of a host cell, upstream of an endogenous gene, designated gene A. The activation domain vector is shown integrated into the genome of the same host cell upstream of an endogenous gene, designated gene B. Both vectors are integrated into the genome of the same host cell. Following integration, each vector is capable of producing a fusion protein containing the binding domain (or activation domain, as the case may be) and the protein encoded by the downstream endogenous gene. If the binding domain fusion protein interacts with the activation domain fusion protein, a protein complex will be formed. This complex is capable of increasing expression of a reporter gene present in the cell.

FIG. 27. Examples of activation vectors useful for in vitro and in vivo transposition. Each vector is illustrated schematically in its linearized form. Each horizontal line represents a DNA molecule. The arrows denote promoter sequences. Boxes indicate exons. Hatched boxes indicate untranslated regions. The solid boxes indicate the transposon signals. It is recognized that there is directionality to the transposon signals, and that the signals are oriented in the configuration suitable for the type of transposition reaction (integration, inversion, or deletion). The following designations were used: splice donor site (S/D), neomycin resistance gene (Neo), dihydrofolate reductase (DHFR), puromycin resistance gene (Puro), poly (A) signal (pA), and the Epstein Barr Virus origin of replication (ori P). It is also recognized that activation exon may encode amino acids in any reading frame (relative to the splice donor site), allowing activation of endogenous genes with different reading frames.

FIG. 28. Schematic illustration depicting integration of an activation vector into a cloned genomic DNA fragment by in vitro transposition. Each

horizontal line represents a DNA molecule. The cloned genomic DNA is in a BAC vector. The single line represents the genomic DNA and the rectangle depicts the BAC vector sequences. The arrows denote promoter sequences located on the DNA molecule, and face in the direction of transcription. Transcribed regions include all sequences located downstream of each promoter. The vector activation exon is depicted as an open box. Exons from a gene encoded in the cloned genomic fragment are depicted as hatched boxes. The solid boxes indicate the transposon signals. It is recognized that there is directionality to the transposon signals, and that the signals are oriented in the configuration suitable for the type of transposition reaction (integration, inversion, or deletion). The following designations were used: splice donor site (S/D), and polyadenylation signal (pA). To integrate the vector into the genomic fragment, the activation vector is incubated with the cloned genomic DNA in the presence of transposase. Following integration of the activation vector into the genomic fragment, the plasmid may be transfected directly into an appropriate eukaryotic host cell to express the gene located downstream of the vector integration site. Alternatively, the BAC plasmid may be transformed into *E. coli* to produce larger quantities of plasmid for transfection into the appropriate eukaryotic host cell.

FIG. 29A-29B. Nucleotide sequence of pRIG14.

FIG. 30A-30C. Nucleotide sequence of pRIG19.

FIG. 31A-31C. Nucleotide sequence of pRIG20.

FIG. 32A-32C. Nucleotide sequence of pRIGad1.

FIG. 33A-33D. Nucleotide sequence of pRIGbd1.

FIG. 34A-34B. Nucleotide sequence of pUniBAC.

FIG. 35A-35B. Nucleotide sequence of pRIG22.

FIG. 36. Schematic diagram of pRIG-TP. The vector is shown in its linearized form. The horizontal line represents a DNA molecule. The arrows denote promoters. Open boxes indicate exons. Filled boxes represent transposon recombination signals (from Tn5 – compatible with the *in vitro* transposition kit available from Epicentre Technologies). The following designations were used:

splice donor site (S/D), puromycin resistance gene (puro), dihydrofolate reductase gene (DHFR), Epstein Barr nuclear antigen - 1 replication protein (EBNA-1), Epstein Barr virus origin of replication (ori P), poly (A) signal (pA), and activation exon (AE). It is understood that the activation exon can contain any sequence capable of directing protein synthesis, including a translation start codon in any reading frame, a partial secretion signal sequence, an entire secretion signal sequence, an epitope tag, a protein, a portion of a protein, or a protein motif. The activation exon may also lack a translation start codon.

FIG. 37A-37C. Nucleotide sequence of pRIG-T.

DETAILED DESCRIPTION OF THE INVENTION

There are great advantages to gene activation by non-homologous recombination over other gene activation procedures. Unlike previous methods of protein over-expression, the methods described herein do not require that the gene of interest be cloned (isolated from the cell). Nor do they require any knowledge of the DNA sequence or structure of the gene to be over-expressed (i.e., the sequence of the ORF, introns, exons, or upstream and downstream regulatory elements) or knowledge of a gene's expression patterns (i.e., tissue specificity, developmental regulation, etc.). Furthermore, the methods do not require any knowledge pertaining to the genomic organization of the gene of interest (i.e., the intron and exon structure).

The methods of the present invention thus involve vector constructs that do not contain target nucleotide sequences for homologous recombination. A target sequence allows homologous recombination of vector DNA with cellular DNA at a predetermined site on the cellular DNA, the site having homology for sequences in the vector, the homologous recombination at the predetermined site resulting in the introduction of the transcriptional regulatory sequence into the genome and the subsequent endogenous gene activation.

The method of the present invention does not involve integration of the vector at predetermined sites. Instead, the present methods involve integration of the vector constructs of the invention into cellular DNA (e.g., the cellular genome) by nonhomologous or "illegitimate" recombination, also called "non-targeted gene activation." In related embodiments, the present invention also concerns non-targeted gene activation. Non-targeted gene activation has a number of important applications. First, by activating genes that are not normally expressed in a given cell type, it becomes possible to isolate a cDNA copy of genes independent of their normal expression pattern. This facilitates isolation of genes that are normally expressed in rare cells, during short developmental periods, and/or at very low levels. Second, by translationally activating genes, it is possible to produce protein expression libraries without the need for cloning the full-length cDNA. These libraries can be screened for new enzymes and proteins and/or for interesting phenotypes resulting from over-expression of an endogenous gene. Third, cell-lines over-expressing a specific protein can be created and used to produce commercial quantities of protein. Thus, activating endogenous genes provides a powerful approach to discovering and isolating new genes and proteins, and to producing large amounts of specific proteins for commercialization.

The vectors described herein do not contain target sequences. A target sequence is a sequence on the vector that has homology with a sequence or sequences within the gene to be activated or upstream of the gene to be activated, the upstream region being up to and including the first functional splice acceptor site on the same coding strand of the gene of interest, and by means of which homology the transcriptional regulatory sequence that activates the gene of interest is integrated into the genome of the cell containing the gene to be activated. In the case of an enhancer integration vector for activating an endogenous gene, the vector does not contain homology to any sequence in the genome upstream or downstream of the gene of interest (or within the gene of interest) for a distance extending as far as enhancer function is operative.

The present methods, therefore, are capable of identifying new genes that have been or can be missed using conventional and currently available cloning techniques. By using the constructs and methodology described herein, unknown and/or uncharacterized genes can be rapidly identified and over-expressed to produce proteins. The proteins have use as, among other things, human therapeutics and diagnostics and as targets for drug discovery.

The methods are also capable of producing over-expression of known and/or characterized genes for *in vitro* or *in vivo* protein production.

A "known" gene is directed to the level of characterization of a gene. The invention allows expression of genes that have been characterized, as well as expression of genes that have not been characterized. Different levels of characterization are possible. These include detailed characterization, such as cloning, DNA, RNA, and/or protein sequencing, and relating the regulation and function of the gene to the cloned sequence (e.g., recognition of promoter and enhancer sequences, functions of the open reading frames, introns, and the like). Characterization can be less detailed, such as having mapped a gene and related function, or having a partial amino acid or nucleotide sequence, or having purified a protein and ascertained a function. Characterization may be minimal, as when a nucleotide or amino acid sequence is known or a protein has been isolated but the function is unknown. Alternatively, a function may be known but the associated protein or nucleotide sequence is not known or is known but has not been correlated to the function. Finally, there may be no characterization in that both the existence of the gene and its function are not known. The invention allows expression of any gene at any of these or other specific degrees of characterization.

Many different proteins (also referred to herein interchangeably as "gene products" or "expression products") can be activated or over-expressed by a single activation construct and in a single set of transfections. Thus, a single cell or different cells in a set of transfectants (library) can over-express more than one protein following transfection with the same or different constructs. Previous

activation methods require a unique construct to be created for each gene to be activated.

Further, many different integration sites adjacent to a single gene can be created and tested simultaneously using a single construct. This allows rapid
5 determination of the optimal genomic location of the activation construct for protein expression.

Using previous methods, the 5' end of the gene of interest had to be extensively characterized with respect to sequence and structure. For each activation construct to be produced, an appropriate targeting sequence had to be isolated. Usually, this must be an isogenic sequence isolated from the same person
10 or laboratory strain of animal as the cells to be activated. In some cases, this DNA may be 50 kb or more from the gene of interest. Thus, production of each targeting construct required an arduous amount of cloning and sequencing of the endogenous gene. However, since sequence and structure information is not
15 required for the methods of the present invention, unknown genes and genes with uncharacterized upstream regions can be activated.

This is made possible using *in situ* gene activation using non-homologous recombination of exogenous DNA sequences with cellular DNA. Methods and compositions (e.g., vector constructs) required to accomplish such *in situ* gene
20 activation using non-homologous recombination are provided by the present invention.

DNA molecules can recombine to redistribute their genetic content by several different and distinct mechanisms, including homologous recombination, site-specific recombination, and non-homologous/illegitimate recombination.
25 Homologous recombination involves recombination between stretches of DNA that are highly similar in sequence. It has been demonstrated that homologous recombination involves pairing between the homologous sequences along their length prior to redistribution of the genetic material. The exact site of crossover can be at any point in the homologous segments. The efficiency of recombination
30 is proportional to the length of homologous targeting sequence (Hope,

Development 113:399 (1991); Reddy *et al.*, *J. Virol.* 65:1507 (1991)), the degree of sequence identity between the two recombining sequences (von Melchner *et al.*, *Genes Dev.* 6:919 (1992)), and the ratio of homologous to non-homologous DNA present on the construct (Letson, *Genetics* 117:759 (1987)).

Site-specific recombination, on the other hand, involves the exchange of genetic material at a predetermined site, designated by specific DNA sequences. In this reaction, a protein recombinase binds to the recombination signal sequences, creates a strand scission, and facilitates DNA strand exchange. *Cre/Lox* recombination is an example of site specific recombination.

Non-homologous/illegitimate recombination, such as that used advantageously by the methods of the present invention, involves the joining (exchange or redistribution) of genetic material that does not share significant sequence homology and does not occur at site-specific recombination sequences. Examples of non-homologous recombination include integration of exogenous DNA into chromosomes at non-homologous sites, chromosomal translocations and deletions, DNA end-joining, double strand break repair of chromosome ends, bridge-breakage fusion, and concatemerization of transfected sequences. In most cases, non-homologous recombination is thought to occur through the joining of "free DNA ends." Free ends are DNA molecules that contain an end capable of being joined to a second DNA end either directly, or following repair or processing. The DNA end may consist of a 5' overhang, 3' overhang, or blunt end.

As used herein, retroviral insertion and other transposition reactions are loosely considered forms of non-homologous recombination. These reactions do not involve the use of homology between the recombining molecules. Furthermore, unlike site-specific recombination, these types of recombination reactions do not occur between discrete sites. Instead, a specific protein/DNA complex is required on only one of the recombination partners (i.e., the retrovirus or transposon), with the second DNA partner (i.e., the cellular genome) usually being relatively non-specific. As a result, these "vectors" do not integrate into the

cellular genome in a targeted fashion, and therefore they can be used to deliver the activation construct according to the present invention.

Vector constructs useful for the methods described herein ideally may contain a transcriptional regulatory sequence that undergoes non-homologous recombination with genomic sequences in a cell to over-express an endogenous gene in that cell. The vector constructs of the invention also lack homologous targeting sequences. That is, they do not contain DNA sequences that target host cell DNA and promote homologous recombination at the target site. Thus, integration of the vector constructs of the present invention into the cellular genome occurs by non-homologous recombination, and can lead to over-expression of a cellular gene via the introduced transcriptional regulatory sequence contained on the integrated vector construct.

The invention is generally directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell. The method does not require previous knowledge of the sequence of the endogenous gene or even of the existence of the gene. Where the sequence of the gene to be activated is known, however, the constructs can be engineered to contain the proper configuration of vector elements (e.g., location of the start codon, addition of codons present in the first exon of the endogenous gene, and the proper reading frame) to achieve maximal overexpression and/or the appropriate protein sequence.

In certain embodiments of the invention, the cell containing the vector may be screened for expression of the gene.

The cell over-expressing the gene can be cultured *in vitro* under conditions favoring the production, by the cell, of desired amounts of the gene product of the endogenous gene that has been activated or whose expression has been increased. If desired, the gene product can then be isolated or purified to use, for example, in protein therapy or drug discovery.

Alternatively, the cell expressing the desired gene product can be allowed to express the gene product *in vivo*.

The vector construct can consist essentially of the transcriptional regulatory sequence.

5 Alternatively, the vector construct can consist essentially of the transcriptional regulatory sequence and one or more amplifiable markers.

10 The invention, therefore, is also directed to methods for over-expressing an endogenous gene in a cell, comprising introducing a vector containing a transcriptional regulatory sequence and an amplifiable marker into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell.

The cell containing the vector is screened for over-expression of the gene.

15 The cell over-expressing the gene is cultured such that amplification of the endogenous gene is obtained. The cell can then be cultured *in vitro* so as to produce desired amounts of the gene product of the amplified endogenous gene that has been activated or whose expression has been increased. The gene product can then be isolated and purified.

Alternatively, following amplification, the cell can be allowed to express the endogenous gene and produce desired amounts of the gene product *in vivo*.

20 The vector construct can consist essentially of the transcriptional regulatory sequence and the splice donor sequence.

25 The invention, therefore, is also directed to methods for over-expressing an endogenous gene in a cell comprising introducing a vector containing a transcriptional regulatory sequence and an unpaired splice donor sequence into the cell, allowing the vector to integrate into the genome of the cell by non-homologous recombination, and allowing over-expression of the endogenous gene in the cell.

The cell containing the vector is screened for expression of the gene.

30 The cell over-expressing the gene can be cultured *in vitro* so as to produce desirable amounts of the gene product of the endogenous gene whose expression

has been activated or increased. The gene product can then be isolated and purified.

Alternatively, the cell can be allowed to express the desired gene product *in vivo*.

5 The vector construct can consist essentially of a transcriptional regulatory sequence operably linked to an unpaired splice donor sequence and also containing an amplifiable marker.

Other activation vectors include constructs with a transcriptional regulatory sequence and an exonic sequence containing a start codon; a
10 transcriptional regulatory sequence and an exonic sequence containing a translational start codon and a secretion signal sequence; constructs with a transcriptional regulatory sequence and an exonic sequence containing a translation start codon, and an epitope tag; constructs containing a transcriptional regulatory sequence and an exonic sequence containing a translational start codon,
15 a signal sequence and an epitope tag; constructs containing a transcriptional regulatory sequence and an exonic sequence with a translation start codon, a signal secretion sequence, an epitope tag, and a sequence-specific protease site. In each of the above constructs, the exon on the construct is located immediately upstream of an unpaired splice donor site.

20 The constructs can also contain a regulatory sequence, a selectable marker lacking a poly(A) signal, an internal ribosome entry site (ires), and an unpaired splice donor site (FIG. 4). A start codon, signal secretion sequence, epitope tag, and/or a protease cleavage site may optionally be included between the ires and the unpaired splice donor sequence. When this construct integrates upstream of
25 a gene, the selectable marker will be efficiently expressed since a poly(A) site will be supplied by the endogenous gene. In addition the downstream gene will also be expressed since the ires will allow protein translation to initiate at the downstream open reading frame (i.e. the endogenous gene). Thus, the message produced by this activation construct will be polycistronic. The advantage of this
30 construct is that integration events that do not occur near genes and in the proper

orientation will not produce a drug resistant colony. The reason for this is that without a poly(A) tail (supplied by the endogenous gene), the neomycin resistance gene will not express efficiently. By reducing the number of nonproductive integration events, the complexity of the library can be reduced without affecting its coverage (the number of genes activated), and this will facilitate the screening process.

In another embodiment of this construct, *cre-lox* recombination sequences can be included between the regulatory sequence and the *neo* start codon and between the ires and the unpaired splice donor site (between the ires and the start codon, if present). Following isolation of cells that have activated the gene of interest, the *neo* gene and ires can be removed by transfecting the cells with a plasmid encoding the *cre* recombinase. This would eliminate the production of the polycistronic message and allow the endogenous gene to be expressed directly from the regulatory sequence on the integrated activation construct. Use of *Cre* recombination to facilitate deletion of genetic elements from mammalian chromosomes has been described (Gu *et al.*, *Science* 265:103 (1994); Sauer, *Meth. Enzymology* 225:890-900 (1993)).

Thus, constructs useful in the methods described herein include, but are not limited to, the following (See also Figures 1-4):

- 1) Construct with a regulatory sequence and an exon lacking a translation start codon.
- 2) Construct with a regulatory sequence and an exon lacking a translation start codon followed by a splice donor site.
- 3) Construct with a regulatory sequence and an exon containing a translation start codon in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 4) Construct with a regulatory sequence and an exon containing a translation start codon in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.

- 5) Construct with a regulatory sequence and an exon containing a translation start codon in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 6) Construct with a regulatory sequence and an exon containing a translation start codon and a signal secretion sequence in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 7) Construct with a regulatory sequence and an exon containing a translation start codon and a signal secretion sequence in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.
- 8) Construct with a regulatory sequence and an exon containing a translation start codon and a signal secretion sequence in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 9) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon and an epitope tag in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 10) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon and an epitope tag in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.
- 11) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon and an epitope tag in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 12) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, and an epitope tag in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 13) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, and an epitope tag in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.

- 14) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, and an epitope tag in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 15) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, an epitope tag, and a sequence specific protease site in reading frame 1 (relative to the splice donor site), followed by an unpaired splice donor site.
- 16) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, an epitope tag, and a sequence specific protease site in reading frame 2 (relative to the splice donor site), followed by an unpaired splice donor site.
- 17) Construct with a regulatory sequence and an exon containing (from 5' to 3') a translation start codon, a signal secretion sequence, an epitope tag, and a sequence specific protease site in reading frame 3 (relative to the splice donor site), followed by an unpaired splice donor site.
- 18) Construct with a regulatory sequence linked to a selectable marker, followed by an internal ribosome entry site, and an unpaired splice donor site.
- 19) Construct 18 in which a cre/lox recombination signal is located between a) the regulatory sequence and the open reading frame of the selectable marker and b) between the ires and the unpaired splice donor site.
- 20) Construct with a regulatory sequence operably linked to an exon containing green fluorescent protein lacking a stop codon, followed by an unpaired splice donor site.

It is to be understood, however, that any vector used in the methods described herein can include one or more (*i.e.*, one, two, three, four, five, or more, and most preferably one or two) amplifiable markers. Accordingly, methods can include a step in which the endogenous gene is amplified. Placement of one or more amplifiable markers on the activation construct results in the juxtaposition

of the gene of interest and the one or more amplifiable markers in the activated cell. Once the activated cell has been isolated, expression can be further increased by selecting for cells containing an increased copy number of the locus containing both the gene of interest and the activation construct. This can be accomplished by selection methods known in the art, for example by culturing cells in selective culture media containing one or more selection agents that are specific for the one or more amplifiable markers contained on the genetic construct or vector.

Following activation of an endogenous gene by nonhomologous integration of any of the vectors described above, the expression of the endogenous gene may be further increased by selecting for increased copies of the amplifiable marker(s) located on the integrated vector. While such an approach may be accomplished using one amplifiable marker on the integrated vector, in an alternative embodiment the invention provides such methods wherein two or more (*i.e.*, two, three, four, five, or more, and most preferably two) amplifiable markers may be included on the vector to facilitate more efficient selection of cells that have amplified the vector and flanking gene of interest. This approach is particularly useful in cells that have a functional endogenous copy of one or more of the amplifiable marker(s) that are contained on the vector, since the selection procedure can result in isolation of cells that have incorrectly amplified the endogenous amplifiable marker(s) rather than the vector-encoded amplifiable marker(s). This approach is also useful to select against cells that develop resistance to the selective agent by mechanisms that do not involve gene amplification. The approach using two or more amplifiable markers is advantageous in these situations because the probability of a cell developing resistance to two or more selective agents (resistance to which is encoded by two or more amplifiable markers) without amplifying the integrated vector and flanking gene of interest is significantly lower than the probability of the cell developing resistance to any single selective agent. Thus, by selecting for two or more vector encoded amplifiable markers, either simultaneously or sequentially, a greater

percentage of cells that are ultimately isolated will contain the amplified vector and gene of interest.

Thus, in another embodiment, the vectors of the invention may contain two or more (*i.e.*, two, three, four, five, or more, and most preferably two) amplifiable markers. This approach allows more efficient amplification of the vector sequences and adjacent gene of interest following activation of expression.

Examples of amplifiable markers that may be used constructing the present vectors include, but are not limited to, dihydrofolate reductase, adenosine deaminase, aspartate transcarbamylase, dihydro-orotase, and carbamyl phosphate synthase.

It is also understood that any of the constructs described herein may contain a eukaryotic viral origin of replication, either in place of, or in conjunction with an amplifiable marker. The presence of the viral origin of replication allows the integrated vector and adjacent endogenous gene to be isolated as an episome and/or amplified to high copy number upon introduction of the appropriate viral replication protein. Examples of useful viral origins include, but are not limited to, SV40 ori and EBV ori P.

The invention also encompasses embodiments in which the constructs disclosed herein consist essentially of the components specifically described for these constructs. It is also understood that the above constructs are examples of constructs useful in the methods described herein, but that the invention encompasses functional equivalents of such constructs.

The term "vector" is understood to generally refer to the vehicle by which the nucleotide sequence is introduced into the cell. It is not intended to be limited to any specific sequence. The vector could itself be the nucleotide sequence that activates the endogenous gene or could contain the sequence that activates the endogenous gene. Thus, the vector could be simply a linear or circular polynucleotide containing essentially only those sequences necessary for activation, or could be these sequences in a larger polynucleotide or other construct such as a DNA or RNA viral genome, a whole virion, or other biological

construct used to introduce the critical nucleotide sequences into a cell. It is also understood that the phrase "vector construct" or the term "construct" may be used interchangeably with the term "vector" herein.

The vector can contain DNA sequences that exist in nature or that have been created by genetic engineering or synthetic processes.

The construct, upon nonhomologous integration into the genome of a cell, can activate expression of an endogenous gene. Expression of the endogenous gene may result in production of full length protein, or in production of a truncated biologically active form of the endogenous protein, depending on the integration site (e.g., upstream region versus intron 2). The activated gene may be a known gene (e.g., previously cloned or characterized) or unknown gene (previously not cloned or characterized). The function of the gene may be known or unknown.

Examples of proteins with known activities include, but are not limited to, cytokines, growth factors, neurotransmitters, enzymes, structural proteins, cell surface receptors, intracellular receptors, hormones, antibodies, and transcription factors. Specific examples of known proteins that can be produced by this method include, but are not limited to, erythropoietin, insulin, growth hormone, glucocerebrosidase, tissue plasminogen activator, granulocyte-colony stimulating factor (G-CSF), granulocyte/macrophage colony stimulating factor (GM-CSF), macrophage colony-stimulating factor (M-CSF) interferon α , interferon β , interferon γ , interleukin-2, interleukin-3, interleukin-4, interleukin-6, interleukin-8, interleukin-10, interleukin-11, interleukin-12, interleukin-13, interleukin-14, TGF- β , blood clotting factor V, blood clotting factor VII, blood clotting factor VIII, blood clotting factor IX, blood clotting factor X, TSH- β , bone growth factor-2, bone growth factor-7, tumor necrosis factor, alpha-1 antitrypsin, anti-thrombin III, leukemia inhibitory factor, glucagon, Protein C, protein kinase C, stem cell factor, follicle stimulating hormone β , urokinase, nerve growth factors, insulin-like growth factors, insulinotropin, parathyroid hormone, lactoferrin, complement inhibitors, platelet derived growth factor, keratinocyte

growth factor, hepatocyte growth factor, endothelial cell growth factor, neurotrophin-3, thrombopoietin, chorionic gonadotropin, thrombomodulin, alpha glucosidase, epidermal growth factor, and fibroblast growth factor. The invention also allows the activation of a variety of genes expressing transmembrane proteins, and production and isolation of such proteins, including but not limited to cell surface receptors for growth factors, hormones, neurotransmitters and cytokines such as those described above, transmembrane ion channels, cholesterol receptors, receptors for lipoproteins (including LDLs and HDLs) and other lipid moieties, integrins and other extracellular matrix receptors, cytoskeletal anchoring proteins, immunoglobulin receptors, CD antigens (including CD2, CD3, CD4, CD8, and CD34 antigens), and other cell surface transmembrane structural and functional proteins that are known in the art. As one of ordinary skill will appreciate, other cellular proteins and receptors that are known in the art may also be produced by the methods of the invention.

One of the advantages of the method described herein is that virtually any gene can be activated. However, since genes have different genomic structures, including different intron/exon boundaries and locations of start codons, a variety of activation constructs is provided to activate the maximum number of different genes within a population of cells.

These constructs can be transfected separately into cells to produce libraries. Each library contains cells with a unique set of activated genes. Some genes will be activated by several different activation constructs. In addition, portions of a gene can be activated to produce truncated, biologically active proteins. Truncated proteins can be produced, for example, by integration of an activation construct into introns or exons in the middle of an endogenous gene rather than upstream of the second exon.

Use of different constructs also allows the activated gene to be modified to contain new sequences. For example, a secretion signal sequence can be included on the activation construct to facilitate the secretion of the activated gene. In some cases, depending on the intron/exon structure or the gene of

interest, the secretion signal sequence can replace all or part of the signal sequence of the endogenous gene. In other cases, the signal sequence will allow a protein which is normally located intracellularly to be secreted.

The regulatory sequence on the vector can be a constitutive promoter. Alternatively, the promoter may be inducible. Use of inducible promoters will allow low basal levels of activated protein to be produced by the cell during routine culturing and expansion. The cells may then be induced to produce large amounts of the desired proteins, for example, during manufacturing or screening. Examples of inducible promoters include, but are not limited to, the tetracycline inducible promoter and the metallothionein promoter.

In preferred embodiments of the invention, the regulatory sequence on the vectors of the invention may be a promoter, an enhancer, or a repressor, any of which may be tissue specific.

The regulatory sequence on the vector can be isolated from cellular or viral genomes. Examples of cellular regulatory sequences include, but are not limited to, regulatory elements from the actin gene, metallothionein I gene, immunoglobulin genes, casein I gene, serum albumin gene, collagen gene, globin genes, laminin gene, spectrin gene, ankyrin gene, sodium/potassium ATPase gene, and tubulin gene. Examples of viral regulatory sequences include, but are not limited to, regulatory elements from *Cytomegalovirus* (CMV) immediate early gene, adenovirus late genes, SV40 genes, retroviral LTRs, and *Herpesvirus* genes. Typically, regulatory sequences contain binding sites for transcription factors such as NF-kB, SP-1, TATA binding protein, AP-1, and CAAT binding protein. Functionally, the regulatory sequence is defined by its ability to promote, enhance, or otherwise alter transcription of an endogenous gene.

In certain preferred embodiments, the regulatory sequence is a viral promoter. In particularly preferred embodiments, the promoter is the CMV immediate early gene promoter. In alternative embodiments, the regulatory element is a cellular, non-viral promoter.

In alternative preferred embodiments, the regulatory element may be or may contain an enhancer. In particularly preferred such embodiments, the enhancer is the cytomegalovirus immediate early gene enhancer. In alternative embodiments, the enhancer is a cellular, non-viral enhancer.

In alternative preferred embodiments, the regulatory element may be or may contain a repressor. In particularly preferred such embodiments, the repressor may be a viral repressor or a cellular, non-viral repressor.

The transcriptional regulatory sequence can also comprise one or more scaffold-attachment regions or matrix attachment sites, negative regulatory elements, and transcription factor binding sites. Regulatory sequences can also include locus control regions.

The invention also encompasses the use of retrovirus transcriptional regulatory sequences, e.g., long terminal repeats. Where these are used, however, they are not necessarily linked to any retrovirus sequence that materially affects the function of the transcriptional regulatory sequence as a promoter or enhancer of transcription of the endogenous gene to be activated (i.e., the cellular gene with which the transcriptional regulatory sequence recombines to activate).

The vector constructs of the invention may also comprise a regulatory sequence which is not operably linked to exonic sequences on the vector. For example, when the regulatory element is an enhancer, it can integrate near an endogenous gene (e.g., upstream, downstream, or in an intron) and stimulate expression of the gene from its endogenous promoter. By this mechanism of activation, exonic sequences from the vector are absent in the transcript of the activated gene.

Alternatively, the regulatory element may be operably linked to an exon. The exon may be a naturally occurring sequence or may be non-naturally occurring (e.g., produced synthetically). To activate endogenous genes lacking a start codon in their first exon (e.g., follicle stimulating hormone- β), a start codon is preferably omitted from the exon on the vector. To activate endogenous genes containing a start codon in the first exon (e.g., erythropoietin and growth

hormone), the exon on the vector preferably contains a start codon, usually ATG and preferably an efficient translation initiation site (Kozak, *J. Mol Biol.* 196: 947 (1987)). The exon may contain additional codons following the start codon. These codons may be derived from a naturally occurring gene or may be non-naturally occurring (e.g., synthetic). The codons may be the same as the codons present in the first exon of the endogenous gene to be activated. Alternatively, the codons may be different than the codons present in the first exon of the endogenous gene. For example, the codons may encode an epitope tag, signal secretion sequence, transmembrane domain, selectable marker, or screenable marker. Optionally, an unpaired splice donor site may be present immediately 3' of the exonic sequence. When the structure of the gene to be activated is known, the splice donor site should be placed adjacent to the vector exon in a location such that the codons in the vector will be in frame with the codons of the second exon of the endogenous gene following splicing. When the structure of the endogenous gene to be activated is not known, separate constructs, each containing a different reading frame, are used.

Operably linked is defined as a configuration that allows transcription through the designated sequence(s). For example, a regulatory sequence that is operably linked to an exonic sequence indicates that the exonic sequence is transcribed. When a start codon is present on the vector, operably linked also indicates that the open reading frame from the vector exon is in frame with the open reading frame of the endogenous gene. Following nonhomologous integration, the regulatory sequence (e.g., a promoter) on the vector becomes operably linked to an endogenous gene and facilitates transcription initiation, at a site generally referred to as a CAP site. Transcription proceeds through the exonic elements on the vector (and, if present, through the start codon, open reading frame, and/or unpaired splice donor site), and through the endogenous gene. The primary transcript produced by this operable linkage is spliced to create a chimeric transcript containing exonic sequences from both the vector and the

endogenous gene. This transcript is capable of producing the endogenous protein when translated.

5 An exon or "exonic sequence" is defined as any transcribed sequence that is present in the mature RNA molecule. The exon on the vector may contain untranslated sequences, for example, a 5' untranslated region. Alternatively, or in conjunction with the untranslated sequences, the exon may contain coding sequences such as a start codon and open reading frame. The open reading frame can encode naturally occurring amino acid sequences or non-naturally occurring amino acid sequences (e.g., synthetic codons). The open reading frame may also
10 encode a signal secretion sequence, epitope tag, exon, selectable marker, screenable marker, or nucleotides that function to allow the open reading frame to be preserved when spliced to an endogenous gene.

Splicing of primary transcripts, the process by which introns are removed, is directed by a splice donor site and a splice acceptor site, located at the 5' and 3' ends of introns, respectively. The consensus sequence for splice donor sites is
15 (A/C)AG GURAGU (where R represents a purine nucleotide) with nucleotides in positions 1-3 located in the exon and nucleotides GURAGU located in the intron.

20 An unpaired splice donor site is defined herein as a splice donor site present on the activation construct without a downstream splice acceptor site. When the vector is integrated by nonhomologous recombination into a host cell's genome, the unpaired splice donor site becomes paired with a splice acceptor site from an endogenous gene. The splice donor site from the vector, in conjunction with the splice acceptor site from the endogenous gene, will then direct the excision of all of the sequences between the vector splice donor site and the
25 endogenous splice acceptor site. Excision of these intervening sequences removes sequences that interfere with translation of the endogenous protein.

The terms upstream and downstream, as used herein, are intended to mean in the 5' or in the 3' direction, respectively, relative to the coding strand. The
30 term "upstream region" of a gene is defined as the nucleic acid sequence 5' of its

second exon (relative to the coding strand) up to and including the last exon of the first adjacent gene having the same coding strand. Functionally, the upstream region is any site 5' of the second exon of an endogenous gene capable of allowing a nonhomologously integrated vector to become operably linked to the endogenous gene.

The vector construct can contain a selectable marker to facilitate the identification and isolation of cells containing a nonhomologously integrated activation construct. Examples of selectable markers include genes encoding neomycin resistance (neo), hypoxanthine phosphoribosyl transferase (HPRT), puromycin (pac), dihydro-otase glutamine synthetase (GS), histidine D (his D), carbamyl phosphate synthase (CAD), dihydrofolate reductase (DHFR), multidrug resistance 1 (mdr1), aspartate transcarbamylase, xanthine-guanine phosphoribosyl transferase (gpt), and adenosine deaminase (ada).

Alternatively, the vector can contain a screenable marker, in place of or in addition to, the selectable marker. A screenable marker allows the cells containing the vector to be isolated without placing them under drug or other selective pressures. Examples of screenable markers include genes encoding cell surface proteins, fluorescent proteins, and enzymes. The vector containing cells may be isolated, for example, by FACS using fluorescently-tagged antibodies to the cell surface protein or substrates that can be converted to fluorescent products by a vector encoded enzyme.

Alternatively, selection can be effected by phenotypic selection for a trait provided by the endogenous gene product. The activation construct, therefore, can lack a selectable marker other than the "marker" provided by the endogenous gene itself. In this embodiment, activated cells can be selected based on a phenotype conferred by the activated gene. Examples of selectable phenotypes include cellular proliferation, growth factor independent growth, colony formation, cellular differentiation (e.g., differentiation into a neuronal cell, muscle cell, epithelial cell, etc.), anchorage independent growth, activation of cellular factors (e.g., kinases, transcription factors, nucleases, etc.), expression of cell

surface receptors/proteins, gain or loss of cell-cell adhesion, migration, and cellular activation (e.g., resting versus activated T cells).

A selectable marker may also be omitted from the construct when transfected cells are screened for gene activation products without selecting for the stable integrants. This is particularly useful when the efficiency of stable integration is high.

The vector may contain one or more (*i.e.*, one, two, three, four, five, or more, and most preferably one or two) amplifiable markers to allow for selection of cells containing increased copies of the integrated vector and the adjacent activated endogenous gene. Examples of amplifiable markers include but are not limited to dihydrofolate reductase (DHFR), adenosine deaminase (*ada*), dihydro-orotase glutamine synthetase (*GS*), and carbamyl phosphate synthase (*CAD*).

The vector may contain eukaryotic viral origins of replication useful for gene amplification. These origins may be present in place of, or in conjunction with, an amplifiable marker.

The vector may also contain genetic elements useful for the propagation of the construct in micro-organisms. Examples of useful genetic elements include microbial origins of replication and antibiotic resistance markers.

These vectors, and any of the vectors disclosed herein, and obvious variants recognized by one of ordinary skill in the art, can be used in any of the methods described herein to form any of the compositions producible by those methods.

Nonhomologous integration of the construct into the genome of a cell results in the operable linkage between the regulatory elements from the vector and the exons from an endogenous gene. In preferred embodiments, the insertion of the vector regulatory sequences is used to upregulate expression of the endogenous gene. Upregulation of gene expression includes converting a transcriptionally silent gene to a transcriptionally active gene. It also includes enhancement of gene expression for genes that are already transcriptionally active,

but produce protein at levels lower than desired. In other embodiments, expression of the endogenous gene may be affected in other ways such as downregulation of expression, creation of an inducible phenotype, or changing the tissue specificity of expression.

5 According to the invention, *in vitro* methods of production of a gene expression product may comprise, for example, (a) introducing a vector of the invention into a cell; (b) allowing the vector to integrate into the genome of the cell by non-homologous recombination; (c) allowing over-expression of an endogenous gene in the cell by upregulation of the gene by the transcriptional regulatory sequence contained on the vector; (d) screening the cell for over-expression of the endogenous gene; and (e) culturing the cell under conditions favoring the production of the expression product of the endogenous gene by the cell. Such *in vitro* methods of the invention may further comprise isolating the expression product to produce an isolated gene expression product. 10 In such methods, any art-known method of protein isolation may be advantageously used, including but not limited to chromatography (e.g., HPLC, FPLC, LC, ion exchange, affinity, size exclusion, and the like), precipitation (e.g., ammonium sulfate precipitation, immunoprecipitation, and the like), electrophoresis, and other methods of protein isolation and purification that will be familiar to one of ordinary skill in the art. 15 20

Analogously, *in vivo* methods of production of a gene expression product may comprise, for example, (a) introducing a vector of the invention into a cell; (b) allowing the vector to integrate into the genome of the cell by non-homologous recombination; (c) allowing over-expression of an endogenous gene in the cell by upregulation of the gene by the transcriptional regulatory sequence contained on the vector; (d) screening the cell for over-expression of the endogenous gene; and (e) introducing the isolated and cloned cell into a eukaryote under conditions favoring the overexpression of the endogenous gene by the cell *in vivo* in the eukaryote. According to this aspect of the invention, any eukaryote 25 30 may be advantageously used, including fungi (particularly yeasts), plants, and

animals, more preferably animals, still more preferably vertebrates, and most preferably mammals, particularly humans. In certain related embodiments, the invention provides such methods which further comprise isolating and cloning the cell prior to introducing it into the eukaryote.

5 As used herein the phrases "conditions favoring the production" of an expression product, "conditions favoring the overexpression" of a gene, and "conditions favoring the activation" of a gene, in a cell or by a cell *in vitro* refer to any and all suitable environmental, physical, nutritional or biochemical parameters that allow, facilitate, or promote production of an expression product, or overexpression or activation of a gene, by a cell *in vitro*. Such conditions may, of course, include the use of culture media, incubation, lighting, humidity, etc., that are optimal or that allow, facilitate, or promote production of an expression product, or overexpression or activation of a gene, by a cell *in vitro*. Analogously, as used herein the phrases "conditions favoring the production" of an expression product, "conditions favoring the overexpression" of a gene, and "conditions favoring the activation" of a gene, in a cell or by a cell *in vivo* refer to any and all suitable environmental, physical, nutritional, biochemical, behavioral, genetic, and emotional parameters under which an animal containing a cell is maintained, that allow, facilitate, or promote production of an expression product, or overexpression or activation of a gene, by a cell in a eukaryote *in vivo*. Whether a given set of conditions are favorable for gene expression, activation, or overexpression, *in vitro* or *in vivo*, may be determined by one of ordinary skill using the screening methods described and exemplified below, or other methods for measuring gene expression, activation, or overexpression that are routine in the art.

25 As used herein, the phrase "activating an endogenous gene" means inducing the production of a transcript encoding the endogenous gene at levels higher than those normally found in the cell containing the endogenous gene. In some applications, "activating an endogenous gene" may also mean producing the

protein, or a portion of the protein, encoded by the endogenous gene at levels higher than those normally found in the cell containing the endogenous gene.

The invention also encompasses cells made by any of the above methods. The invention encompasses cells containing the vector constructs, cells in which the vector constructs have integrated, and cells which are over-expressing desired gene products from an endogenous gene, over-expression being driven by the introduced transcriptional regulatory sequence.

Cells used in this invention can be derived from any eukaryotic species and can be primary, secondary, or immortalized. Furthermore, the cells can be derived from any tissue in the organism. Examples of useful tissues from which cells can be isolated and activated include, but are not limited to, liver, kidney, spleen, bone marrow, thymus, heart, muscle, lung, brain, testes, ovary, islet, intestinal, bone marrow, skin, bone, gall bladder, prostate, bladder, embryos, and the immune and hematopoietic systems. Cell types include fibroblast, epithelial, neuronal, stem, and follicular. However, any cell or cell type can be used to activate gene expression using this invention.

The methods can be carried out in any cell of eukaryotic origin, such as fungal, plant or animal. Preferred embodiments include vertebrates and particularly mammals, and more particularly, humans.

The construct can be integrated into primary, secondary, or immortalized cells. Primary cells are cells that have been isolated from a vertebrate and have not been passaged. Secondary cells are primary cells that have been passaged, but are not immortalized. Immortalized cells are cell lines that can be passaged, apparently indefinitely.

In preferred embodiments, the cells are immortalized cell lines. Examples of immortalized cell lines include, but are not limited to, HT1080, HeLa, Jurkat, 293 cells, KB carcinoma, T84 colonic epithelial cell line, Raji, Hep G2 or Hep 3B hepatoma cell lines, A2058 melanoma, U937 lymphoma, and WI38 fibroblast cell line, somatic cell hybrids, and hybridomas.

Cells used in this invention can be derived from any eukaryotic species, including but not limited to mammalian cells (such as rat, mouse, bovine, porcine, sheep, goat, and human), avian cells, fish cells, amphibian cells, reptilian cells, plant cells, and yeast cells. Preferably, overexpression of an endogenous gene or gene product from a particular species is accomplished by activating gene expression in a cell from that species. For example, to overexpress endogenous human proteins, human cells are used. Similarly, to overexpress endogenous bovine proteins, for example bovine growth hormone, bovine cells are used.

The cells can be derived from any tissue in the eukaryotic organism. Examples of useful vertebrate tissues from which cells can be isolated and activated include, but are not limited to, liver, kidney, spleen, bone marrow, thymus, heart, muscle, lung, brain, immune system (including lymphatic), testes, ovary, islet, intestinal, stomach, bone marrow, skin, bone, gall bladder, prostate, bladder, zygotes, embryos, and hematopoietic tissue. Useful vertebrate cell types include, but are not limited to, fibroblasts, epithelial cells, neuronal cells, germ cells (*i.e.*, spermatocytes/spermatozoa and oocytes), stem cells, and follicular cells. Examples of plant tissues from which cells can be isolated and activated include, but are not limited to, leaf tissue, ovary tissue, stamen tissue, pistil tissue, root tissue, tubers, gametes, seeds, embryos, and the like. One of ordinary skill will appreciate, however, that any eukaryotic cell or cell type can be used to activate gene expression using the present invention.

Any of the cells produced by any of the methods described are useful for screening for expression of a desired gene product and for providing desired amounts of a gene product that is over-expressed in the cell. The cells can be isolated and cloned.

Cells produced by this method can be used to produce protein *in vitro* (e.g., for use as a protein therapeutic) or *in vivo* (e.g., for use in cell therapy).

Commercial growth and production conditions often vary from the conditions used to grow and prepare cells for analytical use (e.g., cloning, protein or nucleic acid sequencing, raising antibodies, X-ray crystallography analysis,

enzymatic analysis, and the like). Scale up of cells for growth in roller bottles involves increase in the surface area on which cells can attach. Microcarrier beads are, therefore, often added to increase the surface area for commercial growth. Scale up of cells in spinner culture may involve large increases in volume. Five liters or greater can be required for both microcarrier and spinner growth. Depending on the inherent potency (specific activity) of the protein of interest, the volume can be as low as 1-10 liters. 10-15 liters is more common. However, up to 50-100 liters may be necessary and volume can be as high as 10,000-15,000 liters. In some cases, higher volumes may be required. Cells can also be grown in large numbers of T flasks, for example 50-100.

Despite growth conditions, protein purification on a commercial scale can also vary considerably from purification for analytic purposes. Protein purification in a commercial practical context can be initially the mass equivalent of 10 liters of cells at approximately 10^4 cells/ml. Cell mass equivalent to begin protein purification can also be as high as 10 liters of cells at up to 10^6 or 10^7 cells/ml. As one of ordinary skill will appreciate, however, a higher or lower initial cell mass equivalent may also be advantageously used in the present methods.

Another commercial growth condition, especially when the ultimate product is used clinically, is cell growth in serum-free medium, by which is intended medium containing no serum or not in amounts that are required for cell growth. This obviously avoids the undesired co-purification of toxic contaminants (e.g., viruses) or other types of contaminants, for example, proteins that would complicate purification. Serum-free media for growth of cells, commercial sources for such media, and methods for cultivation of cells in serum-free media, are well-known to those of ordinary skill in the art.

A single cell made by the methods described above can over-express a single gene or more than one gene. More than one gene can be activated by the integration of a single construct or by the integration of multiple constructs in the same cell (i.e., more than one type of construct). Therefore, a cell can contain

only one type of vector construct or different types of constructs, each capable of activating an endogenous gene.

The invention is also directed to methods for making the cells described above by one or more of the following: introducing one or more of the vector constructs; allowing the introduced construct(s) to integrate into the genome of the cell by non-homologous recombination; allowing over-expression of one or more endogenous genes in the cell; and isolating and cloning the cell.

The term "transfection" has been used herein for convenience when discussing introducing a polynucleotide into a cell. However, it is to be understood that the specific use of this term has been applied to generally refer to the *introduction* of the polynucleotide into a cell and is also intended to refer to the introduction by other methods described herein such as electroporation, liposome-mediated introduction, retrovirus-mediated introduction, and the like (as well as according to its own specific meaning).

The vector can be introduced into the cell by a number of methods known in the art. These include, but are not limited to, electroporation, calcium phosphate precipitation, DEAE dextran, lipofection, and receptor mediated endocytosis, polybrene, particle bombardment, and microinjection. Alternatively, the vector can be delivered to the cell as a viral particle (either replication competent or deficient). Examples of viruses useful for the delivery of nucleic acid include, but are not limited to, adenoviruses, adeno-associated viruses, retroviruses, Herpesviruses, and vaccinia viruses. Other viruses suitable for delivery of nucleic acid molecules into cells that are known to one of ordinary skill may be equivalently used in the present methods.

Following transfection, the cells are cultured under conditions, as known in the art, suitable for nonhomologous integration between the vector and the host cell's genome. Cells containing the nonhomologously integrated vector can be further cultured under conditions, as known in the art, allowing expression of activated endogenous genes.

The vector construct can be introduced into cells on a single DNA construct or on separate constructs and allowed to concatemerize.

Whereas in preferred embodiments, the vector construct is a double-stranded DNA vector construct, vector constructs also include single-stranded DNA, combinations of single- and double-stranded DNA, single-stranded RNA, double-stranded RNA, and combinations of single- and double-stranded RNA. Thus, for example, the vector construct could be single-stranded RNA which is converted to cDNA by reverse transcriptase, the cDNA converted to double-stranded DNA, and the double-stranded DNA ultimately recombining with the host cell genome.

In preferred embodiments, the constructs are linearized prior to introduction into the cell. Linearization of the activation construct creates free DNA ends capable of reacting with chromosomal ends during the integration process. In general, the construct is linearized downstream of the regulatory element (and exonic and splice donor sequences, if present). Linearization can be facilitated by, for example, placing a unique restriction site downstream of the regulatory sequences and treating the construct with the corresponding restriction enzyme prior to transfection. While not required, it is advantageous to place a "spacer" sequence between the linearization site and the proximal most functional element (e.g., the unpaired splice donor site) on the construct. When present, the spacer sequence protects the important functional elements on the vector from exonucleolytic degradation during the transfection process. The spacer can be composed of any nucleotide sequence that does not change the essential functions of the vector as described herein.

Circular constructs can also be used to activate endogenous gene expression. It is known in the art that circular plasmids, upon transfection into cells, can integrate into the host cell genome. Presumably, DNA breaks occur in the circular plasmid during the transfection process, thereby generating free DNA ends capable of joining to chromosome ends. Some of these breaks in the construct will occur in a location that does not destroy essential vector functions

(e.g., the break will occur downstream of the regulatory sequence), and therefore, will allow the construct to be integrated into a chromosome in a configuration capable of activating an endogenous gene. As described above, spacer sequences may be placed on the construct (e.g., downstream of the regulatory sequences).
5 During transfection, breaks that occur in the spacer region will create free ends at a site in the construct suitable for activation of an endogenous gene following integration into the host cell genome.

The invention also encompasses libraries of cells made by the above described methods. A library can encompass all of the clones from a single transfection experiment or a subset of clones from a single transfection experiment. The subset can over-express the same gene or more than one gene, for example, a class of genes. The transfection can have been done with a single type of construct or with more than one type of construct.
10

A library can also be formed by combining all of the recombinant cells from two or more transfection experiments, by combining one or more subsets of cells from a single transfection experiment or by combining subsets of cells from separate transfection experiments. The resulting library can express the same gene, or more than one gene, for example, a class of genes. Again, in each of these individual transfections, a unique construct or more than one construct can be used.
15
20

Libraries can be formed from the same cell type or different cell types.

The library can be composed of a single type of cell containing a single type of activation construct which has been integrated into chromosomes at spontaneous DNA breaks or at breaks generated by radiation, restriction enzymes, and/or DNA breaking agents, applied either together (to the same cells) or separately (applied to individual groups of cells and then combining the cells together to produce the library). The library can be composed of multiple types of cells containing a single or multiple constructs which were integrated into the genome of a cell treated with radiation, restriction enzymes, and/or DNA breaking agents, applied either together (to the same cells) or separately (applied to
25
30

individual groups of cells and then combining the cells together to produce the library).

The invention is also directed to methods for making libraries by selecting various subsets of cells from the same or different transfection experiments. For example, all of the cells expressing nuclear factors (as determined by the presence of nuclear green fluorescent protein in cells transfected with construct 20) can be pooled to create a library of cells with activated nuclear factors. Similarly, cells expressing membrane or secreted proteins can be pooled. Cells can also be grouped by phenotype, for example, growth factor independent growth, growth factor independent proliferation, colony formation, cellular differentiation (e.g., differentiation into a neuronal cell, muscle cell, epithelial cell, etc.), anchorage independent growth, activation of cellular factors (e.g., kinases, transcription factors, nucleases, etc.), gain or loss of cell-cell adhesion, migration, or cellular activation (e.g., resting versus activated T cells).

The invention is also directed to methods of using libraries of cells to over-express an endogenous gene. The library is screened for the expression of the gene and cells are selected that express the desired gene product. The cell can then be used to purify the gene product for subsequent use. Expression of the cell can occur by culturing the cell *in vitro* or by allowing the cell to express the gene *in vivo*.

The invention is also directed to methods of using libraries to identify novel gene and gene products.

The invention is also directed to methods for increasing the efficiency of gene activation by treating the cells with agents that stimulate or effect the patterns of non-homologous integration. It has been demonstrated that gene expression patterns, chromatin structure, and methylation patterns can differ dramatically from cell type to cell type. Even different cell lines from the same cell type can have significant differences. These differences can impact the patterns of non-homologous integration by affecting both the DNA breakage pattern and the repair process. For example, chromatinized stretches of DNA (characteristics

likely associated with inactive genes) may be more resistant to breakage by restriction enzymes and chemical agents, whereas they may be susceptible to breakage by radiation.

Furthermore, inactive genes can be methylated. In this case, restriction enzymes that are blocked by CpG methylation will be unable to cleave methylated sites near the inactive gene, making it more difficult to activate that gene using methylation-sensitive enzymes. These problems can be circumvented by creating activation libraries in several cell lines using a variety of DNA breakage agents. By doing this, a more complete integration pattern can be created and the probability of activating a given gene maximized.

The methods of the invention can include introducing double strand breaks into the DNA of the cell containing the endogenous gene to be over-expressed. These methods introduce double-strand breaks into the genomic DNA in the cell prior to or simultaneously with vector integration. The mechanism of DNA breakage can have a significant effect on the pattern of DNA breaks in the genome. As a result, DNA breaks produced spontaneously or artificially with radiation, restriction enzymes, bleomycin, or other breaking agents, can occur in different locations.

In order to increase integration efficiency and to improve the random distribution of integration sites, cells can be treated with low, intermediate, or high doses of radiation prior to or following transfection. By artificially inducing double strand breaks, the transfected DNA can now integrate into the host cell chromosome as part of the DNA repair process. Normally, creation of double strand breaks to serve as the site of integration is the rate limiting step. Thus, by increasing chromosome breaks using radiation (or other DNA damaging agents), a larger number of integrants can be obtained in a given transfection. Furthermore, the mechanism of DNA breakage by radiation is different than by spontaneous breakage.

Radiation can induce DNA breaks directly when a high energy photon hits the DNA molecule. Alternatively, radiation can activate compounds in the cell

which in turn, react with and break the DNA strand. Spontaneous breaks, on the other hand, are thought to occur by the interaction between reactive compounds produced in the cell (such as superoxides and peroxides) and the DNA molecule. However, DNA in the cell is not present as a naked, deproteinized polymer, but instead is bound to chromatin and present in a condensed state. As a result, some regions are not accessible to agents in the cell that cause double strand breaks. The photons produced by radiation have wave lengths short enough to hit highly condensed regions of DNA, thereby inducing breaks in DNA regions that are under represented in spontaneous breaks. Thus, radiation is capable of creating different DNA breakage patterns, which in turn, should lead to different integration patterns.

As a result, libraries produced using the same activation construct in cells with and without radiation treatment will potentially contain different sets of activated genes. Finally, radiation treatment increases efficiency of nonhomologous integration by up to 5-10 fold, allowing complete libraries to be created using fewer cells. Thus, radiation treatment increases the efficiency of gene activation and generates new integration and activation patterns in transfected cells. Useful types of radiation include α , β , γ , x-ray, and ultraviolet radiation. Useful doses of radiation vary for different cell types, but in general, dose ranges resulting in cell viabilities of 0.1% to >99% are useful. For HT1080 cells, this corresponds to radiation doses from a ^{137}Cs source of approximately 0.1 rads to 1000 rads. Other doses may also be useful as long as the dose either increases the integration frequency or changes the pattern of integration sites.

In addition to radiation, restriction enzymes can be used to artificially induce chromosome breaks in transfected cells. As with radiation, DNA restriction enzymes can create chromosome breaks which, in turn, serve as integration sites for the transfected DNA. This larger number of DNA breaks increases the overall efficiency of integration of the activation construct. Furthermore, the mechanism of breakage by restriction enzymes differs from that by radiation, the pattern of chromosome breaks is also likely to be different.

Restriction enzymes are relatively large molecules compared to photons and small metabolites capable of damaging DNA. As a result, restriction enzymes will tend to break regions that are less condensed than the genome as a whole. If the gene of interest lies within an accessible region of the genome, then treatment of the cells with a restriction enzyme can increase the probability of integrating the activation construct upstream of the gene of interest. Since restriction enzymes recognize specific sequences, and since a given restriction site may not lie upstream of the gene of interest, a variety of restriction enzymes can be used. It may also be important to use a variety of restriction enzymes since each enzyme has different properties (e.g., size, stability, ability to cleave methylated sites, and optimal reaction conditions) that affect which sites in the host chromosome will be cleaved. Each enzyme, due to the different distribution of cleavable restriction sites, will create a different integration pattern.

Therefore, introduction of restriction enzymes (or plasmids capable of expressing restriction enzymes) before, during, or after introduction of the activation construct will result in the activation of different sets of genes. Finally, restriction enzyme-induced breaks increase the integration efficiency by up to 5-10 fold (Yorifuji *et al.*, *Mut. Res.* 243:121 (1990)), allowing fewer cells to be transfected to produce a complete library. Thus, restriction enzymes can be used to create new integration patterns, allowing activation of genes which failed to be activated in libraries produced by non-homologous recombination at spontaneous breaks or at other artificially induced breaks.

Restriction enzymes can also be used to bias integration of the activation construct to a desired site in the genome. For example, several rare restriction enzymes have been described which cleave eukaryotic DNA every 50-1000 kilobases, on average. If a rare restriction recognition sequence happens to be located upstream of a gene of interest, by introducing the restriction enzyme at the time of transfection along with the activation construct, DNA breaks can be preferentially upstream of the gene of interest. These breaks can then serve as sites for integration of the activation construct. Any enzyme can be that cleaves

in an appropriate location in or near the gene of interest and its site is under-represented in the rest of the genome or its site is over-represented near genes (e.g., restriction sites containing CpG). For genes that have not been previously identified, restriction enzymes with 8 bp recognition sites (e.g., *NotI*, *SfiI*, *PmeI*, *SwaI*, *SseI*, *SrfI*, *SgrAI*, *PacI*, *AscI*, *SgfI*, and *Sse8387I*), enzymes recognizing CpG containing sites (e.g., *EagI*, *Bsi-*WI**, *MluI*, and *BssHII*) and other rare cutting enzymes can be used.

In this way, "biased" libraries can be created which are enriched for certain types of activated genes. In this respect, restriction enzyme sites containing CpG dinucleotides are particularly useful since these sites are under-represented in the genome at large, but over-represented in the form of CpG islands at the 5' end of many genes, the very location that is useful for gene activation. Enzymes recognizing these sites, therefore, will preferentially cleave at the 5' end of genic sequences.

Restriction enzymes can be introduced into the host cell by several methods. First, restriction enzymes can be introduced into the cell by electroporation (Yorifuji *et al.*, *Mut. Res.* 243:121 (1990); Winegar *et al.*, *Mut. Res.* 225:49 (1989)). In general, the amount of restriction enzyme introduced into the cell is proportional to its concentration in the electroporation media. The pulse conditions must be optimized for each cell line by adjusting the voltage, capacitance, and resistance. Second, the restriction enzyme can be expressed transiently from a plasmid encoding the enzyme under the control of eukaryotic regulatory elements. The level of enzyme produced can be controlled by using inducible promoters, and varying the strength of induction. In some cases, it may be desirable to limit the amount of restriction enzyme produced (due to its toxicity). In these cases, weak or mutant promoters, splice sites, translation start codons, and poly(A) tails can be utilized to lower the amount of restriction enzyme produced. Third, restriction enzymes can be introduced by agents that fuse with or permeabilize the cell membrane. Liposomes and streptolysin O (Pimplikar *et al.*, *J. Cell Biol.* 125:1025 (1994)) are examples of this type of

agent. Finally, mechanical perforation (Beckers *et al.*, *Cell* 50:523-534 (1987)) and microinjection can also be used to introduce nucleases and other proteins into cells. However, any method capable of delivering active enzymes to a living cell is suitable.

DNA breaks induced by bleomycin and other DNA damaging agents can also produce DNA breakage patterns that are different. Thus, any agent or incubation condition capable of generating double strand breaks in cells is useful for increasing the efficiency and/or altering the sites of non-homologous recombination. Examples of classes of chemical DNA breaking agents include, but are not limited to, peroxides and other free radical generating compounds, alkylating agents, topoisomerase inhibitors, anti-neoplastic drugs, acids, substituted nucleotides, and enediyne antibiotics.

Specific chemical DNA breaking agents include, but are not limited to, bleomycin, hydrogen peroxide, cumene hydroperoxide, tert-butyl hydroperoxide, hypochlorous acid (reacted with aniline, 1-naphthylamine or 1-naphthol), nitric acid, phosphoric acid, doxorubicin, 9-deoxydoxorubicin, demethyl-6-deoxyrubicin, 5-iminodaunorubicin, adriamycin, 4'-(9-acridinylamino)methanesulfonm-aniside, neocarzinostatin, 8-methoxycaffeine, etoposide, ellipticine, iododeoxyuridine, and bromodeoxyuridine.

It has been shown that DNA repair machinery in the cell can be induced by pre-exposing the cell to low doses of a DNA breaking agent such as radiation or bleomycin. By pretreating cells with these agents approximately 24 hours prior to transfection, the cell will be more efficient at repairing DNA breaks and integrating DNA following transfection. In addition, higher doses of radiation or other DNA breaking agents can be used since the LD50 (the dose that results in lethality in 50% of the exposed cells) is higher following pretreatment. This allows random activation libraries to be created at multiple doses and results in a different distribution of integration sites within the host cell's chromosomes.

Screening

Once an activation library (or libraries) is created, it can be screened using a number of assays. Depending on the characteristics of the protein(s) of interest (e.g., secreted versus intracellular proteins) and the nature of the activation construct used to create the library, any or all of the assays described below can be utilized. Other assay formats can also be used.

ELISA. Activated proteins can be detected using the enzyme-linked immunosorbent assay (ELISA). If the activated gene product is secreted, culture supernatants from pools of activation library cells are incubated in wells containing bound antibody specific for the protein of interest. If a cell or group of cells has activated the gene of interest, then the protein will be secreted into the culture media. By screening pools of library clones (the pools can be from 1 to greater than 100,000 library members), pools containing a cell(s) that has activated the gene of interest can be identified. The cell of interest can then be purified away from the other library members by sib selection, limiting dilution, or other techniques known in the art. In addition to secreted proteins, ELISA can be used to screen for cells expressing intracellular and membrane-bound proteins. In these cases, instead of screening culture supernatants, a small number of cells is removed from the library pool (each cell is represented at least 100-1000 times in each pool), lysed, clarified, and added to the antibody-coated wells.

ELISA Spot Assay. ELISA spot are coated with antibodies specific for the protein of interest. Following coating, the wells are blocked with 1% BSA/PBS for 1 hour at 37°C. Following blocking, 100,000 to 500,000 cells from the random activation library are applied to each well (representing ~10% of the total pool). In general, one pool is applied to each well. If the frequency of a cell expressing the protein of interest is 1 in 10,000 (i.e., the pool consists of 10,000 individual clones, one of which expresses the protein of interest), then plating 500,000 cells per well will yield 50 specific cells. Cells are incubated in the wells at 37°C for 24 to 48 hours without being moved or disturbed. At the end of the incubation, the cells are removed and the plate is washed 3 times with PBS/0.05%

Tween 20 and 3 times with PBS/1%BSA. Secondary antibodies are applied to the wells at the appropriate concentration and incubated for 2 hours at room temperature or 16 hours at 4°C. These antibodies can be biotinylated or labeled directly with horseradish peroxidase (HRP). The secondary antibodies are removed and the plate is washed with PBS/1% BSA. The tertiary antibody or streptavidin labeled with HRP is added and incubated for 1 hour at room temperature.

FACS assay. The fluorescence-activated cell sorter (FACS) can be used to screen the random activation library in a number of ways. If the gene of interest encodes a cell surface protein, then fluorescently-labeled antibodies are incubated with cells from the activation library. If the gene of interest encodes a secreted protein, then cells can be biotinylated and incubated with streptavidin conjugated to an antibody specific to the protein of interest (Manz *et al.*, *Proc. Natl. Acad. Sci. (USA)* 92:1921 (1995)). Following incubation, the cells are placed in a high concentration of gelatin (or other polymer such as agarose or methylcellulose) to limit diffusion of the secreted protein. As protein is secreted by the cell, it is captured by the antibody bound to the cell surface. The presence of the protein of interest is then detected by a second antibody which is fluorescently labeled. For both secreted and membrane bound proteins, the cells can then be sorted according to their fluorescence signal. Fluorescent cells can then be isolated, expanded, and further enriched by FACS, limiting dilution, or other cell purification techniques known in the art.

Magnetic Bead Separation. The principle of this technique is similar to FACS. Membrane bound proteins and captured secreted proteins (as described above) are detected by incubating the activation library with an antibody-conjugated magnetic beads that are specific for the protein of interest. If the protein is present on the surface of a cell, the magnetic beads will bind to that cell. Using a magnet, the cells expressing the protein of interest can be purified away from the other cells in the library. The cells are then released from the beads, expanded, analyzed, and further purified if necessary.

RT-PCR. A small number of cells (equivalent to at least the number of individual clones in the pool) is harvested and lysed to allow purification of the RNA. Following isolation, the RNA is reversed-transcribed using reverse transcriptase. PCR is then carried out using primers specific for the cDNA of the gene of interest.

Alternatively, primers can be used that span the synthetic exon in the activation construct and the exon of the endogenous gene. This primer will not hybridize to and amplify the endogenously expressed gene of interest. Conversely, if the activation construct has integrated upstream of the gene of interest and activated gene expression, then this primer, in conjunction with a second primer specific for the gene will amplify the activated gene by virtue of the presence of the synthetic exon spliced onto the exon from the endogenous gene. Thus, this method can be used to detect activated genes in cells that normally express the gene of interest at lower than desired levels.

Phenotypic Section. In this embodiment, cells can be selected based on a phenotype conferred by the activated gene. Examples of phenotypes that can be selected for include proliferation, growth factor independent growth, colony formation, cellular differentiation (e.g., differentiation into a neuronal cell, muscle cell, epithelial cell, etc.), anchorage independent growth, activation of cellular factors (e.g., kinases, transcription factors, nucleases, etc.), gain or loss of cell-cell adhesion, migration, and cellular activation (e.g., resting versus activated T cells). Isolation of activated cells demonstrating a phenotype, such as those described above, is important because the activation of an endogenous gene by the integrated construct is presumably responsible for the observed cellular phenotype. Thus, the activated gene may be an important therapeutic drug or drug target for treating or inducing the observed phenotype.

The sensitivity of each of the above assays can be effectively increased by transiently upregulating gene expression in the library cells. This can be accomplished for NF- κ B site-containing promoters (on the activation construct) by adding PMA and tumor necrosis factor- α , e.g., to the library. Separately, or

in conjunction with PMA and TNF- α , sodium butyrate can be added to further enhance gene expression. Addition of these reagents can increase expression of the protein of interest, thereby allowing a lower sensitivity assay to be used to identify the gene activated cell of interest.

5 Since large activation libraries are created to maximize activation of many genes, it is advantageous to organize the library clones in pools. Each pool can consist of 1 to greater than 100,000 individual clones. Thus, in a given pool, many activated proteins are produced, often in dilute concentrations (due to the overall size of the pool and the limited number of cells within the pool that produce a given activated protein). Thus, concentration of the proteins prior to screening effectively increases the ability to detect the activated proteins in the screening assay. One particularly useful method of concentration is ultrafiltration; however, other methods can also be used. For example, proteins can be concentrated non-specifically, or semi-specifically by adsorption onto ion exchange, hydrophobic, dye, hydroxyapatite, lectin, and other suitable resins under conditions that bind most or all proteins present. The bound proteins can then be removed in a small volume prior to screening. It is advantageous to grow the cells in serum free media to facilitate the concentration of proteins.

10 In another embodiment, a useful sequence that can be included on the activation construct is an epitope tag. The epitope tag can consist of an amino acid sequence that allows affinity purification of the activated protein (e.g., on immunoaffinity or chelating matrices). Thus, by including an epitope tag on the activation construct, all of the activated proteins from an activation library can be purified. By purifying the activated proteins away from other cellular and media proteins, screening for novel proteins and enzyme activities can be facilitated. In some instances, it may be desirable to remove the epitope tag following purification of the activated protein. This can be accomplished by including a protease recognition sequence (e.g., Factor IIa or enterokinase cleavage site) downstream from the epitope tag on the activation construct. Incubation of the

purified, activated protein(s) with the appropriate protease will release the epitope tag from the proteins(s).

In libraries in which an epitope tag sequence is located on the activation construct, all of the activated proteins can be purified away from all other cellular and media proteins using affinity purification. This not only concentrates the activated proteins, but also purifies them away from other activities that can interfere with the assay used to screen the library.

Once a pool of clones containing cells over-expressing the gene of interest is identified, steps can be taken to isolate the activated cell. Isolation of the activated cell can be accomplished by a variety of methods known in the art. Examples of cell purification methods include limiting dilution, fluorescence activated cell sorting, magnetic bead separation, sib selection, and single colony purification using cloning rings.

In preferred embodiments of the invention, the methods include a process wherein the expression product is purified. In highly preferred embodiments, the cells expressing the endogenous gene product are cultured so as to produce amounts of gene product feasible for commercial application, and especially diagnostic and therapeutic and drug discovery uses.

Any vector used in the methods described herein can include an amplifiable marker. Thereby, amplification of both the vector and the DNA of interest (i.e., containing the over-expressed gene) occurs in the cell, and further enhanced expression of the endogenous gene is obtained. Accordingly, methods can include a step in which the endogenous gene is amplified.

Once the activated cell has been isolated, expression can be further increased by amplifying the locus containing both the gene of interest and the activation construct. This can be accomplished by each of the methods described below, either separately or in combination.

Amplifiable markers are genes that can be selected for higher copy number. Examples of amplifiable markers include dihydrofolate reductase, adenosine deaminase, aspartate transcarbamylase, dihydro-orotase, and carbamyl phosphate

synthase. For these examples, the elevated copy number of the amplifiable marker and flanking sequences (including the gene of interest) can be selected for using a drug or toxic metabolite which is acted upon by the amplifiable marker. In general, as the drug or toxic metabolite concentration increases, cells containing fewer copies of the amplifiable marker die, whereas cells containing increased copies of the marker survive and form colonies. These colonies can be isolated, expanded, and analyzed for increased levels of production of the gene of interest.

Placement of an amplifiable marker on the activation construct results in the juxtaposition of the gene of interest and the amplifiable marker in the activated cell. Selection for activated cells containing increased copy number of the amplifiable marker and gene of interest can be achieved by growing the cells in the presence of increasing amounts of selective agent (usually a drug or metabolite). For example, amplification of dihydrofolate reductase (DHFR) can be selected using methotrexate.

As drug-resistant colonies arise at each increasing drug concentration, individual colonies can be selected and characterized for copy number of the amplifiable marker and gene of interest, and analyzed for expression of the gene of interest. Individual colonies with the highest levels of activated gene expression can be selected for further amplification in higher drug concentrations. At the highest drug concentrations, the clones will express greatly increased amounts of the protein of interest.

When amplifying DHFR, it is convenient to plate approximately 1×10^7 cells at several different concentrations of methotrexate. Useful initial concentrations of methotrexate range from approximately 5 nM to 100 nM. However, the optimal concentration of methotrexate must be determined empirically for each cell line and integration site. Following growth in methotrexate containing media, colonies from the highest concentration of methotrexate are picked and analyzed for increased expression of the gene of interest. The clone(s) with the highest concentration of methotrexate are then grown in higher concentrations of methotrexate to select for further amplification

of DHFR and the gene of interest. Methotrexate concentrations in the micromolar and millimolar range can be used for clones containing the highest degree of gene amplification.

Placement of a viral origin of replication(s) (e.g., ori P or SV40 in human cells, and polyoma ori in mouse cells) on the activation construct will result in the juxtaposition of the gene of interest and the viral origin of replication in the activated cell. The origin and flanking sequences can then be amplified by introducing the viral replication protein(s) in trans. For example, when ori P (the origin of replication on Epstein-Barr virus) is utilized, EBNA-1 can be expressed transiently or stably. EBNA-1 will initiate replication from the integrated ori P locus. The replication will extend from the origin bi-directionally. As each replication product is created, it too can initiate replication. As a result, many copies of the viral origin and flanking genomic sequences including the gene of interest are created. This higher copy number allows the cells to produce larger amounts of the gene of interest.

At some frequency, the replication product will recombine to form a circular molecule containing flanking genomic sequences, including the gene of interest. Cells that contain circular molecules with the gene of interest can be isolated by single cell cloning and analysis by Hirt extraction and Southern blotting. Once purified, the cell containing the episomal genomic locus at elevated copy number (typically 10-50 copies) can be propagated in culture. To achieve higher amplification, the episome can be further boosted by including a second origin adjacent to the first in the original construct. For example, T antigen can be used to boost the copy number of ori P/SV40 episomes to a copy number of ~1000 (Heinzel *et al.*, *J. Virol.* 62:3738 (1988)). This substantial increase in copy number can dramatically increase protein expression.

The invention encompasses over-expression of endogenous genes both *in vivo* and *in vitro*. Therefore, the cells could be used *in vitro* to produce desired amounts of a gene product or could be used *in vivo* to provide that gene product in the intact animal.

5 The invention also encompasses the proteins produced by the methods described herein. The proteins can be produced from either known, or previously unknown genes. Examples of known proteins that can be produced by this method include, but are not limited to, erythropoietin, insulin, growth hormone, glucocerebrosidase, tissue plasminogen activator, granulocyte-colony stimulating factor, granulocyte/macrophage colony stimulating factor, interferon α , interferon β , interferon γ , interleukin-2, interleukin-6, interleukin-11, interleukin-12, TGF β , blood clotting factor V, blood clotting factor VII, blood clotting factor VIII, blood clotting factor IX, blood clotting factor X, TSH- β , bone growth factor 2, bone growth factor-7, tumor necrosis factor, alpha-1 antitrypsin, anti-thrombin III, leukemia inhibitory factor, glucagon, Protein C, protein kinase C, macrophage colony stimulating factor, stem cell factor, follicle stimulating hormone β , urokinase, nerve growth factors, insulin-like growth factors, insulinotropin, parathyroid hormone, lactoferrin, complement inhibitors, platelet derived growth factor, keratinocyte growth factor, neurotrophin-3, thrombopoietin, chorionic gonadotropin, thrombomodulin, alpha glucosidase, epidermal growth factor, FGF, macrophage-colony stimulating factor, and cell surface receptors for each of the above-described proteins.

10
15
20 Where the protein product from the activated cell is purified, any method of protein purification known in the art may be employed.

Isolation of Cells Containing Activated Membrane Protein-Encoding Genes

25 Genes that encode membrane associated proteins are particularly interesting from a drug development standpoint. These genes and the proteins they encode can be used, for example, to develop small molecule drugs using combinatorial chemistry libraries and high through-put screening assays. Alternatively, the proteins or soluble forms of the proteins (e.g., truncated proteins lacking the transmembrane region) can be used as therapeutically active agents in humans or animals. Identification of membrane proteins can also be used to identify new ligands (e.g., cytokines, growth factors, and other effector molecules)

using two hybrid approaches or affinity capture techniques. Many other uses of membrane proteins are also possible.

Current approaches to identifying genes that encode integral membrane proteins involve isolation and sequencing of genes from cDNA libraries. Integral membrane proteins are then identified by ORF analysis using hydrophobicity plots capable of identifying the transmembrane region of the protein. Unfortunately, using this approach a gene encoding an integral membrane protein can not be identified unless the gene is expressed in the cells used to produce the cDNA library. Furthermore, many genes are only expressed in very rare cells, during short developmental windows, and/or at very low levels. As a result, these genes can not be efficiently identified using the currently available approaches.

The present invention allows endogenous genes to be activated without any knowledge of the sequence, structure, function, or expression profile of the genes. Using the disclosed methods, genes may be activated at the transcription level only, or at both the transcription and translation levels. As a result, proteins encoded by the activated endogenous gene can be produced in cells containing the integrated vector. Furthermore, using specific vectors disclosed herein, the protein produced from the activated endogenous gene can be modified, for example, to include an epitope tag. Other vectors (e.g., vectors 12-17 described above) may encode a signal peptide followed by an epitope tag. This vector can be used to isolate cells that have activated expression of an integral membrane protein (see Example 5 below). This vector can also be used to direct secretion of proteins that are not normally secreted.

Thus, the invention also is directed to methods for identifying an endogenous gene encoding a cellular integral membrane protein or a transmembrane protein. Such methods of the invention may comprise one or more steps. For example, one such method of the invention may comprise (a) introducing one or more vectors of the invention into a cell; (b) allowing the vector to integrate into the genome of the cell by non-homologous recombination; (c) allowing over-expression of an endogenous gene in the cell by upregulation of

the gene by the transcriptional regulatory sequence contained on the integrated vector construct; (d) screening the cell for over-expression of the endogenous gene; and (e) characterizing the activated gene to determine its identity as a gene encoding a cellular integral membrane protein. In related embodiments, the invention provides such methods further comprising isolating the activated gene from the cell prior to characterizing the activated gene.

To identify genes that encode integral membrane proteins, vectors integrated into the genome of cells will comprise a regulatory sequence linked to an exonic sequence containing a start codon, a signal sequence, and an epitope tag, followed by an unpaired splice donor site. Upon integration and activation of an endogenous gene, a chimeric protein is produced containing the signal peptide and epitope tag from the vector fused to the protein encoded by the downstream exons of the endogenous gene. This chimeric protein, by virtue of the presence of the vector encoded signal peptide, is directed to the secretory pathway where translation of the protein is completed and the protein is secreted. If, however, the activated endogenous gene encodes an integral membrane protein, and the transmembrane region of that gene is encoded by exons located 3' of the vector integration site, then the chimeric protein will go to the cell surface, and the epitope tag will be displayed on the cell surface. Using known methods of cell isolation (for example flow cytometric sorting, magnetic bead cell sorting, immunoadsorption, or other methods that will be familiar to one of ordinary skill in the art), antibodies to the epitope tag can then be used to isolate the cells from the population that display the epitope tag and have activated an integral membrane encoding gene. These cells can then be used to study the function of the membrane protein. Alternatively, the activated gene may then be isolated from these cells using any art-known method, *e.g.*, through hybridization with a DNA probe specific to the vector-encoded exon to screen a cDNA library produced from these cells, or using the genetic constructs described herein.

The epitope tag encoded by the vector exon may be a short peptide capable of binding to an antibody, a short peptide capable of binding to a

substance (e.g., poly histidine/ divalent metal ion supports, maltose binding protein/maltose supports, glutathione S-transferase/glutathione support), or an extracellular domain (lacking a transmembrane domain) from an integral membrane protein for which an antibody or ligand exists. It will be understood, however, that other types of epitope tags that are familiar to one of ordinary skill in the art may be used equivalently in accordance with the invention.

Vectors for Non-targeted Activation of Endogenous Genes

As noted above, non-targeted gene activation has a number of important applications, including activating endogenous genes in host cells which provides a powerful approach to discovering and isolating new genes and proteins, and to producing large amounts of specific proteins for commercialization. For some applications of non-targeted gene activation, it is desirable to create libraries of cells in which each member of the library contains an activation vector integrated into a unique location in the host cell genome, and in which each member of the library has activated a different endogenous gene. Furthermore, it would be desirable to remove cells from the library that contain an integrated vector, but fail to activate an endogenous gene. Since eukaryotic genomes often contain large regions that lack genes, integration of an activation vector into a region devoid of genes can occur frequently. These integrated vectors, however, fail to activate an endogenous gene, and yet are capable of conferring drug resistance on the host cells when a selectable marker (driven by a suitable promoter and followed by a polyadenylation signal) is included on the activation vector. Even more problematic for gene discovery applications, a transcript containing vector sequences is produced in these cells regardless of whether or not a gene has been activated. In cases where a gene has not been activated, these vector sequence-containing transcripts contain non-genic genomic DNA sequences. As a result, when isolating activated genes, one cannot isolate all RNA (or cDNA) molecules that are derived from the integrated vector (i.e. transcripts containing vector sequences), since many of these transcripts do not encode an endogenous gene.

To overcome these difficulties, the present invention provides highly specific vectors and methods that facilitate isolation of vector-activated genes.

These vectors of the invention are useful for activating expression of endogenous genes and for isolating the mRNA and cDNA corresponding to the activated genes. One such vector reduces the number of cells in which the vector integrated into the genome but failed to activate expression from (or transcription through) an endogenous gene. By removing these cells, fewer library members can be created and screened to isolate a given number of activated genes. Furthermore, vector-containing cells that fail to activate gene expression produce an RNA molecule that can interfere with isolation of bona fide activated genes. Thus, the vectors disclosed herein are particularly useful for producing cells suitable for protein over-expression and/or for isolating cDNA molecules corresponding to activated genes. The second type of vector of the invention is useful for isolating exon I from activated endogenous genes. As a result, these vectors can be used to obtain full-length genes from activated RNA transcripts. Each of the functional vector components described herein may be used separately, or in combination with each other.

Poly(A) Trap Activation Vectors

To facilitate isolation of activated genes, the present invention provides novel gene activation vectors that are capable of producing a drug resistant colony, preferentially upon activation of an endogenous gene. Such vectors are referred to herein as "poly(A) trap vectors." Examples of poly(A) trap vectors are shown in Fig. 8A-8F. The nucleotide sequence of one such dual poly(A) trap vector, designated pRIG21b, is shown in Fig. 15A-15B (SEQ ID NO:19). These vectors contain a transcriptional regulatory sequence (which may be any transcriptional regulatory sequence, including but not limited to the promoters, enhancers, and repressors described herein, and which preferably is a promoter or an enhancer, and most preferably a promoter such as a CMV immediate early gene promoter, an SV40 T antigen promoter, a tetracycline-inducible promoter, or a

5 β -actin promoter) operably linked to a selectable marker gene lacking a poly(A) signal. Since the selectable marker gene lacks a polyadenylation signal, its message will not be stable, and the marker gene product will not be efficiently produced. However, if the activation vector integrates upstream of an endogenous gene, the selectable marker can utilize the polyadenylation signal of the endogenous gene, thereby allowing production of the selectable marker protein in sufficient amounts to confer drug resistance. Thus, cells that integrate this activation vector generally form a drug resistant colony only if an endogenous gene has been activated.

10 The poly(A) trap activation vectors can include any selectable or screenable marker. Furthermore, the selectable marker can be expressed from any promoter that is functional in the cells used to create the integration library. Thus, the selectable marker can be expressed by viral or non-viral promoters. Optionally, an unpaired splice donor site may be included in the construct, preferably 3' of the selectable marker to allow the exon encoding the selectable marker to be spliced directly to the exons of the endogenous gene. When a downstream transcriptional regulatory sequence and a splice donor site is included on the vector, the inclusion of a splice donor site adjacent to the selectable marker results in the removal of these downstream elements from the messenger RNA.

20 In a related embodiment, a second transcriptional regulatory sequence (which may be any transcriptional regulatory sequence, including but not limited to the promoters, enhancers, and repressors described herein, and which preferably is a promoter or an enhancer, and most preferably a promoter) may be located downstream of, and in the same orientation as, the selectable marker. Optionally, 25 an unpaired splice donor site may be linked to the downstream transcriptional regulatory sequence. In this configuration, the poly(A) trap vector is capable of producing a message containing the downstream vector-encoded exon spliced to endogenous exons. As described below, these chimeric transcripts can be translated into native or modified protein, depending on the nature of the vector- 30 encoded exon.

As used herein, a "vector-encoded exon" means a region of a vector downstream of the transcriptional regulatory sequence and between the transcription start site and the unpaired splice donor site found on the vector. The vector-encoded exon is present at the 5' end of the transcript containing the endogenous gene in the fully processed message. Analogously, as used herein, a "vector-encoded intron" is the region of the vector located downstream of the unpaired splice donor site. When a linearization site is present on the vector, the vector-encoded intron is the region of the vector that is downstream of the vector-encoded exon between the unpaired splice donor site and the linearization site. The vector-encoded intron is removed from the activated gene transcript during RNA processing.

Splice Acceptor Trap (SAT) Vectors

As an alternative approach for removing cells that fail to activate an endogenous gene, the invention provides additional vectors designated herein as "Splice Acceptor Trap" (SAT) vectors. These vectors are designed to splice from a vector encoded splice donor site to an endogenous splice acceptor. Furthermore, the vectors are designed to produce a product that is toxic to the host cells (or a product that can be selected against) if splicing does not occur. Thus, these vectors facilitate elimination of cells in which the vector-encoded exon failed to splice to an endogenous exon.

The splice acceptor trap vectors can contain both a positive selectable marker and a negative selectable marker gene oriented in the same direction on the vector. As used herein, a positive selectable marker is a gene that, upon expression, produces a protein capable of facilitating the isolation of cells expressing the marker. Analogously, as used herein, a negative selectable marker is a gene that, upon expression, produces a protein capable of facilitating removal of cells expressing the marker.

The positive selectable marker and the negative selectable marker are preferably separated in the vector construct by an unpaired splice donor site. In

other embodiments, however, the positive selectable marker may be fused to the negative selectable marker gene. In this configuration, an unpaired splice donor site is located between the positive and negative selectable marker, such that the reading frame of the negative selectable marker is preserved. The unpaired splice donor site is preferably located at the junction of the positive and negative selectable markers. However, the unpaired splice donor site may be located anywhere in the fusion gene such that upon splicing to an endogenous splice acceptor site, the positive selectable marker will be expressed in an active form and the negative selectable marker will be expressed in an inactive form, or not at all. In this configuration, the positive selectable marker is located upstream of the negative selectable marker.

It will also be apparent to one of ordinary skill in view of the description contained herein that the positive and negative selectable markers on the SAT vector need not be expressed as a fusion protein. In one embodiment, an internal ribosomal entry site (IRES) is inserted between the positive selectable marker and the negative selectable marker. In this configuration, the unpaired splice donor site can be positioned between the two markers, or in the open reading frame of either marker gene such that, upon splicing, the positive selectable marker will be expressed in an active form and the negative selectable marker will be expressed in an inactive form, or not at all. In another embodiment, the positive selectable marker may be driven from a different transcriptional regulatory sequence than the negative selectable marker. In this configuration, the unpaired splice donor site is located in the 5' untranslated region of the negative selectable marker or anywhere in the open reading frame of the negative selectable marker such that, upon splicing, the negative selectable marker will be produced in an inactive form or not at all. Furthermore, when the positive and negative markers are driven from different transcriptional regulatory sequences, the positive selectable marker may be located upstream or downstream of the negative selectable marker, and the positive selectable marker may contain or lack a splice donor site at its 3' end.

The vectors described herein may contain any positive selectable marker. Examples of positive selectable markers useful in this invention include genes encoding neomycin (neo), hypoxanthine phosphoribosyl transferase (HPRT), puromycin (pac), dihydro-orotase, glutamine synthetase (GS), histidine D (his D), carbamyl phosphate synthase (CAD), dihydrofolate reductase (DHFR), multidrug resistance 1 (mdr1), aspartate transcarbamylase, xanthine-guanine phosphoribosyl transferase (gpt), and adenosine deaminase (ada). Alternatively, the vectors may contain a screenable marker in place of the positive selectable marker. Screenable markers include any protein capable of producing a recognizable phenotype in the host cell. Examples of screenable markers included cell surface epitopes (such as CD2) and enzymes (such as β -galactosidase).

The vectors described herein may also, or alternatively, contain any negative selectable marker that can be selected against. Examples of negative selectable markers include hypoxanthine phosphoribosyl transferase (HPRT), thymidine kinase (TK), and diphtheria toxin. The negative selectable marker can also be a screenable marker, such as a cell surface protein or an enzyme. Cells expressing the negative screenable marker may be removed by, for example, Fluorescence Activated Cell Sorting (FACS) or magnetic bead cell sorting.

To isolate cells that have activated expression of an endogenous gene, the cells containing the integrated vector can be placed under the appropriate drug selection. Selection for the positive selectable marker and against the negative selectable marker can occur simultaneously. In another embodiment, selection can occur sequentially. When selection occurs sequentially, selection for the positive selectable marker can occur first, followed by selection against the negative selectable marker. Alternatively, selection against the negative selectable marker can occur first, followed by selection for the positive selectable marker.

The positive and negative markers are expressed by a transcriptional regulatory element located upstream of the translation start site of each gene. When a positive/negative marker fusion gene or an ires sequence is used, a single transcriptional regulatory element drives expression of both markers. A poly(A)

signal may be placed 3' of each selectable marker. If a positive/negative fusion gene is used a single poly(A) signal is positioned 3' of the markers. Alternatively, a poly(A) signal may be excluded from the vector to provide additional specificity for a gene activation event (see dual poly(A)/splice acceptor trap below).

Dual Poly(A)/Splice Acceptor Trap Vectors

To further reduce the number of cells that lack a gene activation event, the invention also provides vectors that confers host cell survival only if the vector-encoded exon has spliced to an exon from an endogenous gene and has acquired a poly(A) signal. These vectors are designated herein as "dual poly(A)/splice acceptor trap vectors" or as "dual poly(A)/SAT vectors." By requiring both splicing and polyadenylation to occur for cell survival, cells that fail to activate an endogenous gene are more efficiently eliminated from the activation library.

The dual poly(A)/splice acceptor trap vectors contain a positive selectable marker and a negative selectable marker configured as described for the SAT vectors; however, neither gene contains a functional poly(A) signal. Thus, the positive selectable marker will not be expressed at high levels unless splicing occurs to capture an endogenous poly(A) signal. Aside from the lack of a poly(A) signal, all other features and embodiments of this type of vector are the same as those of the SAT vectors as described herein. Examples of dual poly(A)/SAT vectors are shown in Figs. 9A-9F and 10A-10F. The nucleotide sequence of one such dual poly(A)/SAT vector, designated pRIG22b, is shown in Fig. 16A-16B (SEQ ID NO:20).

Vectors for Activating Protein Expression from Endogenous Genes

In many applications of non-targeted gene activation, it is desirable to produce protein from the activated endogenous gene. To accomplish this, a second transcriptional regulatory sequence (which may be any transcriptional regulatory sequence, including but not limited to the promoters, enhancers, and repressors described herein, and which is preferably a promoter or an enhancer,

and most preferably a promoter) can be placed downstream of the selectable marker(s) on any of the vectors described herein. When poly(A) trap vectors, SAT vectors, or dual poly(A) trap/SAT vectors are used, the downstream transcriptional regulatory sequence is positioned to drive expression in the same direction as the upstream selectable marker(s). To activate expression of full-length protein with this type of vector, however, the vector must integrate into the 5' UTR of the endogenous gene to avoid cryptic start ATG codons upstream of exon I.

Alternatively, to increase the frequency of protein expression using non-targeted gene activation, the downstream transcriptional regulatory sequence on the vector may be operably linked to an exonic sequence followed by a splice donor site. In a preferred embodiment, the vector exon lacks a start codon. This vector is particularly useful for activating protein expression from genes that do not encode the translation start codon in exon I. In an alternative preferred embodiment, the vector exon contains a start codon. Additional codons can be located between the translational start codon and the splice donor site. For example, a partial signal secretion sequence can be encoded on the vector exon. The partial signal sequence can be any amino acid sequence capable of complementing a partial signal sequence from an endogenous gene to produce a functional signal sequence. The partial sequence may encode between one and one hundred amino acids, and may be derived from existing genes, or may consist of novel sequences. Thus, this vector is useful for producing and secreting protein from genes that encode part of the endogenous signal sequence in exon I, and the remainder in subsequent exons. In another example of a vector useful for activating a particular type of endogenous gene, a functional signal sequence can be encoded on the vector exon. This vector allows protein to be produced and secreted from genes that encode a signal sequence in exon I. It can also be used to produce secreted forms of proteins that are not normally secreted.

In cases where a start codon is included on the vector exon, it can be advantageous to produce a vector in each reading frame. This is achieved by

varying the number of nucleotides between the start codon and the splice donor junction site. Together, the preferred vector configurations are capable of producing protein from endogenous genes, regardless of the exon/intron structure, location of the translation start codon, or reading frame.

5 **Vectors for Isolating Exon I from Activated Endogenous Genes**

10 The non-targeted gene activation vectors described above are useful for activating and isolating endogenous genes and for producing protein from endogenous genes. Upon integration upstream of an endogenous gene, however each of these vectors produces a transcript that lacks exon I from the endogenous gene. Since the vectors are designed to produce a transcript containing the vector encoded exon spliced to the first splice acceptor site downstream of the vector integration site, and since the first exon of eukaryotic genes does not contain a splice acceptor site, normally, the first exon of endogenous genes will not be recovered on mRNA molecules derived from non-targeted gene activation. For
15 some genes, such as genes that contain coding information in the first exon, there is a need to efficiently recover the first exon of the activated endogenous gene.

20 To recover the first exon of activated endogenous genes, a transcriptional regulatory sequence (which may be any transcriptional regulatory sequence, including but not limited to the promoters, enhancers, and repressors described herein, and which is preferably a promoter or an enhancer, and most preferably a promoter) is included on the activation vector downstream of a second transcriptional regulatory sequence (which may also be any transcriptional regulatory sequence, including but not limited to the promoters, enhancers, and repressors described herein, and which is preferably a promoter or an enhancer,
25 and most preferably a promoter) which drives expression of a vector encoded exon. Thus, the upstream transcriptional regulatory sequence is linked to an unpaired splice donor site and the downstream transcriptional regulatory sequence is not linked to a splice donor site. Both transcriptional regulatory sequences are oriented to drive expression in the same direction. Examples of such exon I

recovery vectors are shown in Fig. 12A-12G. The integration of this type of vector will create at least two different types of RNA transcripts (Figure 13). The first transcript is derived from the upstream transcriptional regulatory sequence and contains the vector exon spliced to exon II of an endogenous gene. The second transcript is derived from the downstream transcriptional regulatory sequence and contains, from 5' to 3', the region between the vector and the transcription start site of the gene, exon I, exon II, and all downstream exons. Using methods described herein, both transcripts can be recovered and analyzed, allowing the characterization of exon I from genes isolated by non-targeted gene activation.

The exon located on the activation vector can encode a selectable marker, a protein, a portion of a protein, secretion signal sequences, a portion of a signal sequence, an epitope, or nothing. When a protein is encoded by the exon, a poly(A) signal may be included downstream of the vector encoded gene. Alternatively, a poly(A) signal may be omitted. In another embodiment, a positive and negative selectable marker may be operably linked to the upstream transcriptional regulatory sequence(s). In this embodiment, the position of the unpaired splice donor site relative to the selectable markers is described above for the SAT vectors and the dual poly(A)/SAT vectors.

Gene Activation Vectors for Single-Exon and Multi-Exon Gene Trapping

As noted above, in one embodiment the poly(A) trap vectors of the invention may contain a promoter operably linked to a selectable marker followed by an unpaired splice donor site. Such vectors, when integrated into or near a gene, produce transcripts containing the selectable marker spliced onto an endogenous gene. Since the endogenous gene encodes a poly(A) signal, the resulting mRNA is polyadenylated, thereby allowing the transcript to be translated at levels sufficient to confer drug resistance on the cell containing the integrated vector.

While the vectors described above are capable of "trapping" endogenous genes, the splice donor site downstream of a selectable marker cannot be used in, and in some cases can interfere with, several potential applications for such vectors. First, these vectors cannot be used to selectively trap single exon genes, since these genes do not contain a splice acceptor site. Second, these vectors often "trap" cryptic genes, since drug resistance relies solely on vector integration upstream of a poly (A) signal. Unfortunately, cryptic poly (A) signals exist in the genome, leading to formation of drug resistant cells and creation of non-genic transcripts containing the selectable marker. These cells and transcripts can interfere with gene discovery applications using these vectors. Third, without novel modifications such as those described herein (see above), these vectors are not capable of efficiently producing protein from the activated endogenous gene. Furthermore, protein expression from an endogenous gene can be poor even when an internal ribosome entry site (IRES) is included between the selectable marker and the splice donor site, since translation from an IRES is generally less efficient than translation from the first start codon at the 5' end of a transcript. Thus, there is a need for vectors that are capable of more specifically trapping endogenous genes, including single exon genes, and that are capable of efficiently expressing protein from the activated endogenous genes.

Thus, in additional embodiments, the present invention provides such vectors. In one such embodiment, the vector may contain a promoter operably linked to one or more (*i.e.*, one, two, three, four, five, or more) selectable markers, wherein the selectable marker is not followed by a splice donor site or a poly(A) signal (see Figures 17A-17G). In general, upon integration into a host cell genome, this vector will fail to produce sufficient quantities of selectable marker since the marker transcript will not be polyadenylated. However, if the vector integrates in close proximity to, or into, a gene, including a single exon gene, the selectable marker will acquire a poly(A) signal from the endogenous gene, thereby stabilizing the marker transcript and conferring a drug resistant phenotype on the cell. In addition to selecting for vector integration into or near

genes, vectors according to this aspect of the invention can also be used to recover exon I from the activated gene, as described in the section of this application entitled "Vectors for Isolating Exon I from Activated Endogenous Genes."

In a preferred embodiment, the vector can contain a second selectable marker upstream of the first selectable marker (see Figure 18). The upstream selectable marker is preferably operably linked to a transcriptional regulatory sequence, most preferably a promoter. Optionally, an unpaired splice donor site can be positioned between the transcription start site and the translation start site of the upstream selectable marker. Alternatively, the splice donor site may be located anywhere in the open reading frame of the upstream selectable marker, such that, following vector integration into a host cell genome, and upon splicing from the vector encoded splice donor site to an endogenous exon, the upstream selectable marker will be produced in an inactive form, or not at all. By selecting for cells that produce the downstream positive selectable marker in an active form, cells containing the vector integrated into or near a gene can be isolated. Furthermore, by selecting against cells producing the upstream selectable marker in the active form, cells in which the vector transcript has spliced to an exon from a multi-exon endogenous gene can be removed. In other words, these vectors can be used to isolate cells that contain a vector integrated into a single exon gene or into the 3' most exon of a multi-exon gene since, in these instances, a splice acceptor site is absent between the vector encoded splice donor site and the endogenous poly (A) signal. Thus, the majority of cells containing activated multi-exon genes will not survive selection, and as a result, cells containing activated single exon genes will be greatly enriched in the library.

In another preferred embodiment, vectors according to this aspect of the invention may contain one or more (*i.e.*, one, two, three, four, five, or more, and preferably one) negative selectable marker(s) upstream of the first selectable marker (see Figures 19A and 19B). The negative selectable marker preferably is operably linked to a promoter. Optionally, an unpaired splice donor site may be positioned between the transcription start site and the translation start site of the

negative selectable marker. Alternatively, the splice donor site may be located anywhere in the open reading frame of the negative selectable marker, such that, following vector integration into a host cell genome, and upon splicing from the vector encoded splice donor site to an endogenous exon, the negative selectable marker will be produced in an inactive form, or not at all. By selecting for cells that produce the positive selectable marker in an active form and selecting against cells producing the negative selectable marker in the active form, these vectors can be used to identify cells containing the vector integrated into or upstream of an endogenous gene. Since (1) splicing to an endogenous exon and (2) acquisition of a poly (A) signal are both required for cell survival, cells containing cryptic gene trap events are reduced within the library. The reason for this is that the probability of a vector integrating next to both a cryptic splice acceptor site and a cryptic poly (A) signal is substantially less than the probability of the vector integrating next to a single cryptic site. Thus, these vectors provide a higher degree of specificity for trapping genes than previous vectors.

It will also be recognized by one of ordinary skill in view of the teachings contained herein that vectors containing positive and negative selectable markers can be used to produce protein from the activated endogenous gene. One vector configuration capable of directing protein production consists of the splice donor site positioned in the 5' UTR of the negative selectable marker. Upon splicing, a chimeric transcript containing the 5' UTR from the negative selectable marker linked to the second exon of an endogenous gene is produced. This vector is capable of activating protein production from genes that encode a translation start codon in the second or subsequent exon. Likewise, the splice donor site can be placed in the open reading frame of the negative selectable marker, in a position that does not interfere with the function of the marker unless splicing has occurred. Similar vectors containing the splice donor site positioned in different reading frames relative to the translation start codon can also be used. Upon splicing to an endogenous gene, these vectors will produce a chimeric transcript containing a start codon from the negative selectable marker fused to exon II of

the activated endogenous gene. Thus, these vectors will be capable of activating protein expression from genes that encode a translation start codon in exon I. Additional positive/negative selection vector designs capable of efficiently producing protein from activated endogenous genes are described below.

Any of the vectors of the invention can contain an internal ribosome entry site (ires) 3' of the downstream selectable marker. The ires allows translation of the endogenous gene upon vector integration into an endogenous gene. Optionally, a translation start codon may be included between the selectable marker and the ires sequence. When a start codon is present, additional codons may be present on the exon. The start codon, and if present additional codons, may be present in any, and collectively all, reading frames relative to the splice donor site. Furthermore, the codons downstream of the translation start codon, if present, may encode, for example, a signal secretion signal, a partial signal sequence, a protein (including a full-length protein, a portion of a protein, a protein motif, an epitope tag, etc.), or a spacer region.

In additional preferred embodiments, any of the vectors described herein may contain, upstream of the selectable marker(s), a second transcriptional regulatory sequence (most preferably a promoter) operably linked to a exonic region, followed by an unpaired splice donor site. This upstream exon is particularly useful for expressing protein from activated endogenous genes. The exon may lack a translation start codon. Alternatively, the exon may contain a translation start codon. When a start codon is present, additional codons may be present on the exon. The start codon, and if present additional codons, may be present in any, and collectively all, reading frames relative to the splice donor site. Furthermore, the codons downstream of the translation start codon, if present, may encode, for example, a signal secretion sequence, a partial signal sequence, a protein (including a full-length protein, a portion of a protein, a protein motif, an epitope tag, etc.), or a spacer region.

Activation Vectors Useful for Detecting Protein-protein Interactions

Genetic approaches for detecting protein-protein interactions have previously been described (*see, e.g.*, U.S. Patent Nos. 5,283,173; 5,468,614; and 5,667,973, the disclosures of which are fully incorporated herein by reference). This approach relies on cloning a first cDNA molecule next to, and in frame with, a gene fragment encoding a DNA binding domain; and cloning a second cDNA molecule next to, and in frame with, a gene fragment encoding a transcription transactivation domain. Each chimeric gene is expressed from a promoter region located upstream of the chimeric gene. To detect expression, both chimeric genes are transfected into a reporter cell. If the first chimeric protein interacts with the second chimeric protein (via the proteins encoded by the cloned cDNA's fused to the DNA binding and transcription activation domains), then the DNA binding domain and the transcription activation domain will be joined within a single protein complex. As a result, the protein-protein interaction complex can bind to the regulatory region of the reporter gene and activate its expression.

A limitation of this previous approach is that it is only capable of detecting protein-protein interactions between genes that have been cloned as cDNA. As described herein, many genes are expressed at very low levels, in rare cell types, or during short developmental windows; and therefore, these genes are typically absent from cDNA libraries. Furthermore, many genes are too large to be isolated efficiently as full-length clones, thereby making it difficult to use these previous approaches.

The present invention is capable of activating protein expression from endogenous genes or from transfected genomic DNA. Unlike previous approaches, virtually any gene can be efficiently expressed, regardless of its normal expression pattern. Furthermore, since the present invention is also capable of modifying the protein expressed from the endogenous gene (or from the transfected genomic DNA), it is also possible to produce chimeric proteins for use in protein-protein interaction assays.

To detect protein-protein interactions by the present invention, two vectors are used. The first vector, generally referred to as BD/SD (binding domain/splice donor), contains a promoter operably linked to a polynucleotide encoding a DNA binding domain and an unpaired splice donor site. The second vector, generally referred to as AD/SD (activation domain/splice donor), contains a promoter operably linked to a polynucleotide encoding a transcription activation domain and an unpaired splice donor site. To accommodate genes that have different reading frames, the binding domain and activation domain can be encoded in each of the three possible reading frames relative to the unpaired splice donor site. In addition, BD/SD and AD/SD vectors can have other functional elements, as described herein for other vectors, including selectable markers and amplifiable markers. The vectors may also contain selectable markers oriented in a configuration that permits selection for cells in which the vector has activated a gene. Multi-promoter/activation exon vectors are also useful. Several examples of BD/SD and AD/SD vectors are illustrated in Figure 25. An example illustrating detection of a protein-protein interaction using these vectors is depicted in Figure 26.

The DNA binding domain of the BD/SD vector may encode any protein domain capable of binding to a specific nucleotide sequence. When a transcription activation protein is used to supply the DNA binding domain, the transcription activation domain is omitted from the BD/SD vector. Examples of genes encoding proteins with DNA binding domains include, but are not limited to, the yeast GAL4 gene, the yeast GCN4 gene, and the yeast ADR1 gene. Other genes from prokaryotic and eukaryotic sources may also be used to supply DNA binding domains.

The transcription activation domain of the AD/SD vector encodes a protein domain capable of enhancing transcription of a reporter gene when positioned near the promoter region of the reporter gene. When a transcription activation protein is used to supply the transcription activation domain, the DNA binding domain is omitted from the AD/SD vector. Examples of genes encoding

proteins with transcription activation domains include, but are not limited to, the yeast GAL4 gene, the yeast GCN4 gene, and the yeast ADR1 gene. Other genes from prokaryotic and eukaryotic sources may also be used to supply transcription activation domains.

In the present invention, protein-protein interactions are detected using the BD/SD and AD/SD vectors, described above, to activate expression of genes located in stretches of genomic DNA.

In one embodiment, the BD/SD vector is integrated randomly into the genome of a reporter cell line. As with other vectors described herein, the BD/SD vectors are capable of activating protein expression from genes located downstream of the vector integration site. Since the activation exon on the BD/SD vector encodes a DNA binding domain, the activated endogenous protein will be produced as a fusion protein containing the DNA binding domain at its N-terminus. Thus, by integrating the BD/SD vector into the genome of a host cell, a library of fusion proteins can be created, wherein each protein will contain a DNA binding domain at its N-terminus.

It is also recognized that the AD/SD vector can be integrated into the genome of a reporter cell line to produce a library of cells, wherein each member of the library is expressed as a different endogenous gene fused to a transcription activation domain.

Once created, the BD/SD library may be transfected with a vector expressing a specific gene (referred to below as gene X) fused to a transcription activation domain. This allows virtually any gene encoded in the genome to be tested for an interaction to gene X. Likewise, the AD/SD library may be transfected with a vector expressing a specific gene (e.g. gene X) fused to a DNA binding domain. This allows virtually any gene encoded in the genome to be tested for an interaction to gene X. It is also recognized that the specific gene may be stably expressed in the host cell prior to construction of the BD/SD or AD/SD libraries.

In an alternative embodiment, genomic DNA is cloned into the BD/SD and/or AD/SD vector(s) downstream of the DNA binding domain and activation domain, respectively. If a gene is present and correctly oriented in the genomic DNA, then the BD/SD vector (or the AD/SD vector) will be capable of expressing the gene as a fusion protein useful for detecting protein-protein interactions. Like integration of BD/SD (or AD/SD) vectors *in situ*, any gene can be tested regardless of whether it has been previously isolated as a cDNA molecule.

In another embodiment, a second library is created in the cells of the first library. For example, the AD/SD vector can be integrated into cells comprising the BD/SD library. Conversely, the BD/SD vector can be integrated into cells comprising the AD/SD library. This allows all proteins expressed as binding domain fusion proteins to be tested against all activation domain fusion protein. Since the present invention is capable of expressing substantially all of the proteins (as fusions with the binding and activation domains) in a eukaryotic organism, this approach, for the first time, allows all combinations of protein-protein interactions to be tested in a single library. To survey all protein-protein interactions in an organism, the library within a library must be substantially comprehensive. For example, to detect ~50% of protein-protein interactions in an organism containing 100,000 genes, the first library must contain at least 100,000 cells, each expressing an activated gene. Within each clone of the first library, the second vector would then be used to create a library of at least 100,000 clones, each containing an activated gene. Thus, the total library would contain 100,000 clones x 100,000 clones, or 10^{10} total clones. This assumes all genes are activated at equal frequencies, and that each gene activation event results in production of a fusion protein in frame with the activated endogenous gene. To produce libraries with greater than 50% coverage of protein-protein interactions, and/or to ensure that proteins that are activated at lower frequencies are represented, larger libraries can be created.

It is also recognized that library vs. library screens can be created in several ways. First, both libraries are produced, simultaneously or sequentially, by

integrating BD/SD and AD/SD vectors into the genome of the same reporter cells. Second, a first library is created by integrating a BD/SD vector into the genome of a reporter cell, and a second library is produced by transfecting the AD/SD vector containing cloned genomic DNA. It is recognized that in this approach, the AD/SD library may be created first, followed by introduction of a BD/SD vector containing cloned genomic DNA. It is also recognized that the first library can be created by transfecting the BD/SD vector (or AD/SD vector) containing cloned genomic DNA, followed by integrating the second vector into the reporter cell genome. Third, both libraries are created, simultaneously or sequentially, by transfecting cells with a BD/SD and AD/SD vectors, wherein each vector contains a cloned fragment of genomic DNA. Fourth, it is recognized that when cloned genomic fragments are used in either the BD/SD vector or the AD/SD vector, a cDNA library may be created in the other vector and introduced into cells. This allows all of the genes present in the cDNA library to be tested for interaction with all other genes in the genome.

Since library/library screens involve the creation of large libraries of cells, it is important to maximize the frequency of gene activation and in frame fusion protein production among the members of the library. This can be accomplished in at least two ways. First, the BD/SD and AD/SD vectors can contain selectable markers in a configuration that "traps" genes. Examples of selection trap vectors are shown in Figures 8, 9, 10, 17, 19, 21, and 25. These vectors select for cells in which the activation vector has transcriptionally activated a gene. Second, multiple promoter/activation exon units can be included on the BD/SD and AD/SD vectors. Each promoter/activation exon unit encodes the binding domain (or activation domain) in a different reading frame relative to the unpaired splice donor site. An example of a multi-promoter/exon vector is illustrated in figure 23. This type of vector ensures that any gene activated at the transcription level will be produced as an in frame fusion protein from on of the promoter/activation exon units on the vector. Third, the vectors can be introduced into the reporter cells

using efficient transfection procedures. In this respect, insertion of BD/SD and AD/SD vectors by retroviral integration is advantageous.

Reporter cells useful in the present invention include any cell that is capable of properly splicing the transcripts produced by the BD/SD and AD/SD vectors. The reporter cells contain a reporter gene that is expressed at higher levels in the presence of a protein-protein interaction between proteins expressed from BD/SD and AD/SD vectors. The reporter gene may be a selectable marker, such as any of the markers described herein. Alternatively, the reporter gene may be a screenable marker. Examples of useful selectable markers and screenable markers are described herein.

In the reporter cells, a minimal promoter is operably linked to the reporter gene. To allow increased expression of the reporter gene in the presence of a protein-protein interaction, a DNA binding site is positioned in or near the minimal promoter, such that the DNA binding site is recognized by the protein encoded by the DNA binding domain region of the BD/SD vector. In the absence of a protein-protein interaction, the DNA binding domain fusion protein produced from BD/SD lacks a transcription activation domain, and therefore, can not activate transcription from the minimal promoter of the reporter gene. If, however, the DNA binding domain fusion protein produced from BD/SD interacts with the activation domain fusion protein produced from the AD/SD vector, then the protein complex can activate expression of the reporter gene. Increased reporter gene expression can be detected using an assay for the screenable marker, or using drug selection for a selectable marker.

It is also recognized that other reporter systems can be used in conjunction with the present invention to detect protein-protein interactions. Specifically, any protein that contains two separable domains, each required to be in close proximity with the other to produce a biochemical or structural activity, can be used in conjunction with the present invention.

Multi-Promoter/Activation Exon Vectors

In applications of nontargeted gene activation in which the goal is to activate protein expression from an unknown gene, a collection of vectors typically must be used. Thus, in an additional embodiment, the invention provides vectors containing one or more promoter/activation exon units (see Figures 20A-20E).

To accommodate the variety of gene structures that exist in the genomes of eukaryotic cells, vectors according to this aspect of the invention preferably contain a transcriptional regulatory sequence (*e.g.*, a promoter) operably linked to an activation exon with a different structure. Collectively, these activation exons are capable of activating protein expression from substantially all endogenous genes. For example, to activate protein expression from genes that encode a translation start codon in exon II (or exons downstream of exon II), one vector can contain a transcriptional regulatory sequence (*e.g.*, a promoter) operably linked to an activation exon lacking a translation start codon. To activate protein expression from all types of genes that encode a translation start codon in exon I, three separate vectors must be used, each containing a transcriptional regulatory sequence (*e.g.*, a promoter) operably linked to a different activation exon. Each activation exon encodes a start codon in a different reading frame. Additional activation exon configurations are also useful. For example, to activate protein expression and secretion from genes that encode a portion of their signal secretion sequence in exon I, three separate vectors must be used, each containing a transcriptional regulatory sequence (*e.g.*, a promoter) operably linked to a different activation exon. Each activation exon encodes a partial signal sequence in a different reading frame. To activate protein expression and secretion from genes that encode their entire signal sequence in exon I, three vectors must be used, each containing a transcriptional regulatory sequence (*e.g.*, a promoter) operably linked to a different activation exon. Each activation exon contains an entire signal secretion sequence in a different reading frame. In addition to activating expression of genes that encode secreted proteins,

promoter/activation exons encoding entire signal sequences will also activate expression and secretion of proteins that are not normally secreted. This, for example, can facilitate protein purification of proteins that are normally intracellularly localized.

5 Other useful coding sequences can be included on the activation exon of vectors according to this aspect of the invention, including but not limited to sequences encoding proteins (including full length proteins, portions of proteins, protein motifs, and/or epitope tags). As described herein, vectors according to this aspect of the invention can be integrated, individually or collectively, into the genome of a host cell to produce a library of cells. Each member of the library will potentially overexpress a different endogenous protein. Thus, these collections of vectors make it possible to activate all or substantially all of the endogenous genes in a eukaryotic host cell.

10 When integrating a collection of vectors into host cells, as described above, activation of protein expression can be achieved from substantially any gene. Unfortunately, to produce protein from all endogenous genes, a large number of library members must be generated. In part, this is due to the large number of genes encoded by the host cell. In addition, using this approach, many cells will contain a vector integrated into or near an endogenous gene; however, the integrated vector will contain an activation exon with a structure that is incompatible with activating protein expression from the endogenous gene. For example, the vector exon may encode a start codon in reading frame 1 (relative to the splice junction), whereas the protein encoded by the first exon downstream of the integrated vector may be in reading frame 2 (relative to the splice junction). Thus, many library members will contain an integrated vector that has activated transcription of an endogenous gene, but that failed to produce the protein encoded by the endogenous gene.

25 To decrease the number of cells that fail to activate protein expression following vector integration into or near an endogenous gene, a vector containing multiple promoter/activation exons can be used. On this vector, each

30

promoter/activation exon unit can be capable of activating protein expression from an endogenous gene with a different structure. Since a single vector comprising multiple activation exons is capable of producing multiple transcripts, each containing a different activation exon, a single vector integrated into or near a gene can be capable of activating protein expression, regardless of the structure of the endogenous gene (see Figure 21).

Multi-promoter/activation exon vectors can contain two or more promoter/activation exons. Each promoter/activation exon unit may be followed by an unpaired splice donor site. In one such embodiment, two promoter/activation exons are included on the vector, wherein each promoter/activation exon is capable of activating protein expression from a different type of endogenous gene. In a preferred embodiment, the vector may contain three promoter/activation exons, wherein each exon encodes a translation start codon in a different reading frame. In another preferred embodiment, the vector may contain three promoter/activation exons, wherein each exon encodes a partial signal secretion sequence in a different reading frame. In yet another preferred embodiment, the vector may contain three promoter/activation exons, wherein each exon encodes an entire signal secretion sequence in a different reading frame. Additional embodiments include each of the vectors above containing a fourth promoter/activation exon, wherein the fourth activation exon does not encode a translation start codon.

Any number (*e.g.*, one or more, two or more, three or more, four or more, five or more, etc.) of promoter/activation exon units may be included on the vector. When multiple promoter/activation exons are present on a single vector, they are preferably oriented in the same direction relative to one another (*i.e.*, the promoters drive expression in the same direction).

The promoters that drive transcription of different activation exons may be the same as one another or one or more promoters may be different. The promoters may be viral, cellular, or synthetic. The promoters may be constitutive or inducible. Other types of promoters and regulatory sequences, recognizable to

one skilled in the art or as described herein, may also be used in preparing the vectors according to this aspect of the invention.

Any of the vectors containing multiple promoter/activation exon units may optionally include one or more selectable marker(s) and/or amplifiable marker(s). The selectable and/or amplifiable markers may contain a poly(A) signal. Alternatively, the markers may lack a poly(A) signal. The selectable marker may be a positive or negative selectable marker. The selectable marker may contain an unpaired splice donor site upstream, within, or downstream of the marker. Alternatively, the selectable marker may lack an unpaired splice donor site. The selectable marker(s) and/or amplifiable marker(s), when present, may be located upstream, among, or downstream of the promoter/activation exon units. The selectable and/or amplifiable marker(s) may be located on the vector in any orientation relative to the promoter/activation exon units. When the purpose of the selectable marker is to trap endogenous genes, the selectable marker is preferably oriented in the same direction as the promoter/activation exons.

Amplifiable Markers

Any of the vectors described herein may also optionally comprise one or more (e.g., two, three, four, five, or more) amplifiable markers. Examples of amplifiable markers include those described in detail hereinabove. Preferably, the amplifiable marker(s) are located upstream of the positive/negative selectable marker(s). When using polyadenylation trap vectors, it may be advantageous to omit a polyadenylation signal from the amplifiable marker(s) to eliminate the possibility of capturing a vector-encoded poly(A) signal derived from vector concatemerization prior to integration.

When present, the amplifiable marker(s) may be located upstream of the activation transcriptional regulatory sequence (i.e. the promoter responsible for directing transcription from the vector through the endogenous gene). The amplifiable marker(s) may be present on the vector in any orientation (i.e. the open reading frame may be present on either DNA strand).

It is also understood that the amplifiable marker(s) can also be the same gene as the positive selectable marker. Examples of genes that can be used both as positive selectable markers and amplifiable markers include dihydrofolate reductase, adenosine deaminase (ada), dihydro-*orotase*, glutamine synthase (GS), and carbamyl phosphate synthase (CAD).

In some embodiments and for certain applications, it may be desirable to place multiple amplifiable markers on the vector. Use of more than one amplifiable marker allows dual selection, or alternatively sequential selection, for each amplifiable marker. This facilitates the isolation of cells that have amplified the vector and flanking genomic locus, including the gene of interest.

Promoters

It is understood that any promoter and regulatory element may be used on these activation vectors to drive expression of the selectable marker, amplifiable marker (if present), and/or the endogenous gene. In additional preferred embodiments, the promoter driving expression of the endogenous gene is a strong promoter. The CMV immediate early gene promoter, SV40 T antigen promoter, and β -actin promoter are examples of this type of promoter. In another preferred embodiment, an inducible promoter is used to drive expression of the endogenous genes. This allows endogenous proteins to be expressed in a more controlled fashion. The Tetracycline inducible promoter, heat shock promoter, *ectdyson*e promoter, and metallothionein promoter are examples of this type of promoter. In yet another embodiment, a tissue specific promoter is used to drive expression of endogenous genes. Examples of tissue specific promoters include, but are not limited to, immunoglobulin promoters, casein promoter, and growth hormone promoter.

Restriction Sites

The vectors of the invention can contain one or more restriction sites located downstream of the unpaired splice donor site in the vector. These

restriction sites can be used to linearize plasmid vectors prior to transfection. In the linear configuration, the activation vector contains, from 5' to 3' relative to the transcribed strand, a promoter, a splice donor site, and a linearization site.

5 A restriction site(s) may also be included in the vector intron to facilitate removal of vector intron-containing cDNA molecules. In this embodiment, the vector contains, from 5' to 3' relative to the transcribed strand, a promoter, a splice donor site, a restriction site, and a linearization site. By including a restriction site between the unpaired splice donor site and the linearization site, unspliced transcripts can be removed by digestion of cDNA with the appropriate restriction enzyme. cDNA molecules derived from gene activation have removed the vector intron containing the restriction site, and therefore, will not be digested. This allows gene activated transcripts to be preferentially enriched during amplification/cloning, and greatly facilitates identification and analysis of endogenous genes.

10 A restriction site(s) may also be included in the vector exon to facilitate cloning of activated genes. Following gene activation, mRNA is recovered from cells and synthesized into cDNA. By digesting the cDNA with a restriction enzyme that cuts in the vector exon, gene activated cDNA molecules will contain an appropriate overhang at the 5' end for subsequent cloning into a suitable vector. This facilitates isolation of gene activated cDNA molecules.

15 In one embodiment, the restriction site located in the vector exon is different than the restriction site(s) located in the vector intron. This facilitates removal of cDNA molecules that contain a vector intron since the digested cDNA fragments from vector intron containing transcripts can be designed to have an overhang that is incompatible with the cloning vector (see below). Alternatively, degenerate restriction sites recognized by the same enzyme may be located in the vector exon and intron. Enzymes that cleave these sites are capable of cleaving multiple sites, sites with an odd number of bases in the recognition sequence, sites with interrupted palindromes, nonpalindromic sequences, or sites containing one
20
25
30 or more degenerate bases. In other words, restriction sites recognized by the same

restriction endonuclease may be used if the enzyme produces an overhang in the vector exon that is different from the overhang produced in the vector intron. Since different overhangs are produced, a cloning vector containing a site that is compatible with the vector exon overhang, and incompatible with the vector intron overhang may be used to preferentially clone vector exon containing and vector intron lacking cDNA molecules. Examples of useful degenerate restriction sites include DNA sequences recognized by Sfi I, Acci, Afl III, SapI, Ple I, Tsp45 I, ScrF I, Tse I, PpuM I, Rsr II, and SgrA I.

The restriction site(s) located in the vector intron and/or exon can be a rare restriction site (e.g. an 8 bp restriction site) or an ultra-rare site (e.g. a site recognized by intron encoded nucleases). Examples of restriction enzymes with 8 bp recognitions sites include *NofI*, *SfiI*, *PacI*, *AscI*, *FseI*, *PmeI*, *SgfI*, *SrfI*, *SbfI*, *Sse 8387 I*, and *SwaI*. Examples of intron encoded restriction enzymes include *I-PpoI*, *I-SceI*, *I-CeuI*, *PI-PspI*, and *PI-TliI*. Alternatively, restriction sites smaller than 8 bp can be placed on the vector. For example, restriction sites composed of 7 bp, 6 bp, 5 bp, or 4 bp can be used. In general, the use of smaller the restriction recognition sites will lead to the cloning of less than full-length genes. In some cases, such as creation of hybridization probes, isolation of smaller cDNA clones may be advantageous.

Bidirectional Activation Vectors

The activation vectors described herein can also be bidirectional. When a single activation transcriptional regulatory sequence is present on the vector, gene activation occurs only when the vector integrates into an appropriate location (e.g. upstream of the gene) and in the correct orientation. That is, in order to activate an endogenous gene, the promoter on the activation construct must face the endogenous gene allowing transcription of the coding strand. As a result of this directionality requirement, only half of the integration events into a locus may result in the transcriptional activation of an endogenous gene. The other half of integration events result in the vector transcribing away from a gene of interest.

Therefore, to increase the gene activation frequency by a factor of two, the present invention provides bidirectional vectors that may be used to activate an endogenous gene regardless of the orientation in which the vector integrates into the host cell genome.

5 A bidirectional vector according to this aspect of the invention preferably comprises two transcriptional regulatory sequences (which may be any transcriptional regulatory sequences, including but not limited to the promoters, enhancers, and repressors described herein, and which preferably are promoters or enhancers, and most preferably promoters), two splice donor sites, and a linearization site. When a splice donor site is useful, each transcriptional regulatory sequence is operably linked to a separate splice donor site, and the transcriptional regulatory sequence/splice donor pairs may be in inverse orientation relative to each other (*i.e.*, the first transcriptional regulatory sequence may be integrated into the host cell genome in an orientation that is inverse relative to the orientation in which the second transcriptional regulatory sequence has integrated into the host cell genome). The two opposing transcriptional regulatory sequence/splice donor sites can be separated by the linearization site. The function of the linearization site is to produce free DNA ends between the transcriptional regulatory sequence/splice donor sites (*i.e.* in a location suitable for activation of endogenous genes). Examples of bidirectional vectors of the invention are shown in Fig. 11A-11C.

20 The two opposing transcriptional regulatory sequences may be the same transcriptional regulatory sequences or different transcriptional regulatory sequences. Optionally, a translational start codon (e.g. ATG) and one or more additional codons may be included on either or both vector encoded exons. When a translational start codon is present, either or both vector exons may encode a protein, a portion of a protein, a signal secretion sequence, a portion of a signal secretion sequence, a protein motif, or an epitope tag. Alternatively, either or both vector exons may lack a translational start codon.

The bidirectional vectors according to this aspect of the invention may optionally include one or more selectable markers and one or more amplifiable markers, including those selectable markers and amplifiable markers described in detail herein. The bidirectional vectors may also be configured as poly(A) trap, splice acceptor trap, or dual poly(A)/splice acceptor trap vectors, as described above. Other vector configurations described for unidirectional vectors may also be incorporated into bidirectional vectors.

Co-transfection of Genomic DNA with Non-targeted Activation Vectors

It is recognized that any of the vectors described herein can be integrated into, or otherwise combined with, genomic DNA prior to transfection into a eukaryotic host cell. This permits high level expression from virtually any gene in the genome, regardless of the normal expression characteristics of the gene. Thus, the vectors of the invention can be used to activate expression from genes encoded by isolated genomic DNA fragments. To accomplish this, the vector is integrated into, or otherwise combined with, genomic DNA containing at least one gene, or portion of a gene. Typically, the activation vector must be positioned within or upstream of a gene in order to activate gene expression. Once inserted (or joined), the downstream gene may be expressed (as a transcript or a protein) by introducing the vector/genomic DNA into an appropriate eukaryotic host cell. Following introduction into the host cell, the vector encoded promoter drives expression through the gene encoded in the isolated DNA, and following splicing, produces a mature mRNA molecule. Using appropriate activation vectors, this process allows protein to be expressed from any gene encoded by the transfected genomic DNA. In addition, using the methods described herein, cDNA molecules, corresponding to genes encoded by the transfected genomic DNA, can be generated and isolated.

To achieve stable expression of the activated gene, the transfected activation vector/genomic DNA can be integrated into the host cell genome. Alternatively, the transfected activation vector/genomic DNA can be maintained

as a stable episome (e.g. using a viral origin of replication and/or nuclear retention function - see below). In yet another embodiment, the activated gene may be expressed transiently, for example, from a plasmid.

As used herein, the term "genomic DNA" refers to the unspliced genetic material from a cell. Splicing refers to the process of removing introns from genes following transcription. Thus, genomic DNA, in contrast to mRNA and cDNA, contains exons and introns in an unspliced form. In the present invention, genomic DNA derived from eukaryotic cells is particularly useful since most eukaryotic genes contain exons and introns, and since many of the vectors of the present invention are designed to activate genes encoded in the genomic DNA by splicing to the first downstream exon, and removing intervening introns.

Genomic DNA useful in the present invention may be isolated using any method known in the art. A number of methods for isolating high molecular weight genomic DNA and ultra-high molecular weight genomic DNA (intact and encased in agarose plugs) have been described (Sambrook et al., Molecular Cloning, Cold Spring Harbor Laboratory Press, (1989)). In addition, commercial kits for isolating genomic DNA of various sizes are also available (Gibco/BRL, Stratagene, Clontech, etc.).

The genomic DNA used in the invention may encompass the entire genome of an organism. Alternatively, the genomic DNA may include only a portion of the entire genome from an organism. For example, the genomic DNA may contain multiple chromosomes, a single chromosome, a portion of a chromosome, a genetic locus, a single gene, or a portion of a gene.

Genomic DNA useful in the invention may be substantially intact (i.e. unfragmented) prior to introduction into a host cell. Alternatively, the genomic DNA may be fragmented prior to introduction into a host cell. This can be accomplished by, for example, mechanical shearing, nuclease treatment, chemical treatment, irradiation, or other methods known in the art. When the genomic DNA is fragmented, the fragmentation conditions may be adjusted to produce DNA

fragments of any desirable size. Typically, DNA fragments should be large enough to contain at least one gene, or a portion of a gene (e.g. at least one exon). The genomic DNA may be introduced directly into an appropriate eukaryotic host cell without prior cloning. Alternatively, the genomic DNA (or genomic DNA fragments) may be cloned into a vector prior to transfection. Useful vectors include, but are not limited to, high and intermediate copy number plasmids (e.g. pUC, pBluescript, pACYC184, pBR322, etc.), cosmids, bacterial artificial chromosomes (BAC's), yeast artificial chromosomes (YAC's), P1 artificial chromosomes (PAC's), and phage (e.g. lambda, M13, etc.). Other cloning vectors known in the art may also be used. When genomic DNA has been cloned into a cloning vector, specific cloned DNA fragments may be isolated and used in the present invention. For example, YAC, BAC, PAC, or cosmid libraries can be screened by hybridization to identify clones that map to specific chromosomal regions. Optionally, once isolated, these clones can be ordered to produce a contig through the chromosomal region of interest. To rapidly isolate cDNA copies of the genes present in this contig, these genomic clones may be transfected, separately or en masse, with the activation vector into a host cell. cDNA containing a vector encoded exon, and lacking a vector encoded intron, can then be isolated and analyzed. Thus, since all genes present in a contig can be rapidly isolated as cDNA clones, this approach greatly enhances the speed of positional cloning approaches.

Any activation vector described herein, including derivatives recognized by those skilled in the art, may be co-transfected with genomic DNA, and therefore, are useful in the present invention. In its simplest form, the vector can contain a promoter operably linked to an exon followed by an unpaired splice donor site. Examples of other useful vectors include, but are not limited to, poly A trap vectors (e.g. vectors illustrated in Figures 8, 9, 11C, 12F, and 17), dual poly (A)/Splice acceptor trap vectors (e.g. vectors illustrated in Figures 9, 10, 12G, 19, and 21), bi-directional vectors (e.g. vectors illustrated in Figure 11), single exon trap vectors (e.g. the vector illustrated in Figure 19), multi-

promoter/activation exon vectors (e.g. the vector illustrated in Figure 23), vectors for isolating cDNA's corresponding to activated genes, and vectors for activating protein expression from activated genes (e.g. vectors illustrated in Figures 2, 3, 4, 8B-F, 9B-C, 9E-F, 10B-C, 10E-F, 11, 12, 17B-G, and 23).

5 The activation vector may also contain a viral origin of replication. The presence of a viral origin of replication allows vectors containing genomic fragments to be propagated as an episome in the host cell. Examples of useful viral origins of replication include ori P (Epstein Barr Virus), SV40 ori, BPV ori, and vaccinia ori. To facilitate replication from these origins, the appropriate viral replication proteins may be expressed from the vector. For example, EBV ori P and SV40 ori containing vectors may also encode and express EBNA-1 or T antigen, respectively. Alternatively, the vectors may be introduced into cells that are already expressing the viral replication protein (e.g. EBNA-1 or T antigen). Examples of cells expressing EBNA-1 and T antigen include human 293 cells transfected with an EBNA-1 expression unit (Clontech) and COS-7 cells (American Type Culture Collection; ATCC No. CRL-1651), respectively.

10 The activation vector may also contain an amplifiable marker. This enables cells containing increased copies of the vector and flanking genomic DNA, either episomal or integrated in the host cell genome, to be isolated. Cells containing increased copies of the vector and flanking genomic DNA express the activated gene at higher levels, facilitating gene isolation and protein production.

20 The activation vector and genomic DNA may be introduced into any host cell capable of splicing from the vector-encoded splice donor site to a splice acceptor site encoded by the genomic DNA. In a preferred embodiment, the genomic DNA/activation vector are transfected into a host cell from the same species as the cell from which the genomic DNA was isolated. In some instances, however, it is advantageous to transfect the genomic DNA into a host cell from a species that is different from the cell from which the genomic DNA was isolated. For example, transfection of genomic DNA from one species into a host cell of a second species can facilitate analysis of the genes activated in the transfected

30

genomic DNA using hybridization techniques. Under high stringency hybridization, activated genes that were encoded by the transfected DNA can be distinguished from genes derived from the host cell. Transfection of genomic DNA from one species into a host cell from another species can also be used to produce protein in a heterologous cell. This may allow protein to be produced in heterologous cells that provide growth, protein modification, or manufacturing advantages.

The activation vector may be co-transfected into a host cell along with genomic DNA, wherein the vector is not attached to the genomic DNA prior to introduction into the cell. In this embodiment, the genomic DNA will become fragmented during the transfection process, thereby creating free DNA ends. These DNA ends can become joined to the co-transfected activation vector by the cell's DNA repair machinery. Following joining to the activation vector, the genomic DNA and activation vector can be integrated into the host cell genome by the process of non-homologous recombination. If, during this process, a vector becomes joined to a gene encoded by the transfected genomic DNA, the vector will activate its expression.

Alternatively, the non-targeted activation vector may be physically linked to the genomic DNA prior to transfection. In a preferred embodiment, genomic DNA fragments are ligated to the vector prior to transfection. This is advantageous because it maximizes the probability of the vector becoming operably linked to a gene encoded by the genomic DNA, and minimizes the probability of the vector integrating into the host cell genome without the heterologous genomic DNA.

In a related embodiment, the genomic DNA may be cloned into the activation vector, downstream of the activation exon. In this embodiment, cloning of large genomic fragments can be facilitated in vectors capable of accommodating large genomic fragments. Thus, the activation vector may be constructed in BAC's, YAC's, PAC's, cosmids, or similar vectors capable of propagating large fragments of genomic DNA.

Another method for joining the activation vector to genomic DNA involves transposition. In this embodiment, the activation vector is integrated into the genomic DNA by transposition or retroviral integration reactions prior to transfection into a cell. Accordingly, activation vectors can contain cis sequences necessary for facilitating transposition and/or retroviral integration. Examples of vectors containing transposon signals are illustrated in figure 27; however, it is recognized that any vector described herein may contain transposon signals.

Any transposition system capable of inserting foreign sequences into genomic DNA can be used in the present invention. In addition, transposons capable of facilitating inversions and deletions can also be used to practice the invention. While deletion and inversion systems do not integrate the activation vector into genomic DNA, they do allow the activation vector to change positions relative to cloned genomic DNA when the genomic DNA has been cloned into the activation vector. Thus, multiple genes within a given genomic fragment can be activated by shuffling the activation vector (by integration, inversion, or deletion) into multiple positions within, or outside of, the genomic fragment. Examples of transposition systems useful for the present invention include, but are not limited to $\delta\gamma$, Tn 3, Tn5, Tn7, Tn9, Tn10, Ty, retroviral integration and retro-transposons (Berg et al., *Mobile DNA*, ASM Press, Washington DC, pp. 879-925 (1989); Strathman et al., *Proc. Natl. Acad. Sci. USA* 88:1247 (1991); Berg et al., *Gene* 113:9 (1992); Liu et al., *Nucl. Acids Res.* 15:9461 (1987), Martin et al., *Proc. Natl. Acad. Sci. USA* 92:8398 (1995); Phadnis et al., *Proc. Natl. Acad. Sci. USA* 86:5908 (1989); Tomcsanyi et al., *J. Bacteriol.* 172:6348 (1990); Way et al., *Gene* 32:369 (1984); Bainton et al., *Cell* 65:805 (1991); Ahmed et al., *J. Mol. Biol.* 178:941 (1984); Benjamin et al., *Cell* 59:373 (1989); Brown et al., *Cell* 49:347 (1987); Eichinger et al., *Cell* 54:955 (1988); Eichinger et al., *Genes Dev.* 4:324 (1990); Braiterman et al., *Mol. Cell. Biol.* 14:5719 (1994); Braiterman et al., *Mol. Cell. Biol.* 14:5731 (1994); York et al., *Nucl. Acids Res.* 26:1927 (1998); Devine et al., *Nucl. Acids Res.* 18:3765 (1994); Goryshin et al., *J. Biol. Chem.* 273:7367 (1998)

Using transposition, an activation vector may be integrated into any form of genomic DNA. For example, the activation vector may be integrated into either intact or fragmented genomic DNA. Alternatively, the activation vector may be integrated into a cloned fragment of genomic DNA (Figure 28). In this embodiment, the genomic DNA may reside in any cloning vector, including high and intermediate copy number plasmids (e.g. pUC, pBluescript, pACYC184, pBR322, etc.), cosmids, bacterial artificial chromosomes (BAC's), yeast artificial chromosomes (YAC's), P1 artificial chromosomes (PAC's), and phage (e.g. lambda, M13, etc.). Other cloning vectors known in the art may also be used. As described above, genomic fragments from specific genetic loci may be isolated and used as a substrate for activation vector integration.

Following integration of the activation vector, the genomic DNA may be introduced directly into a suitable host cell for expression of the activated gene. Alternatively, the genomic DNA may be introduced into and propagated in an intermediate host cell. For example, following integration of an activation vector into a BAC genomic library, the BAC library can be transformed into *E. coli*. This allows plasmids containing the transposon to be enriched by selecting for an antibiotic resistance marker residing on the activation vector. As a result, BAC plasmids lacking an integrated activation vector will be removed by antibiotic selection.

The transposition mediated activation vector integration may occur *in vitro* using purified enzymes. Alternatively, the transposition reaction may occur *in vivo*. For example, transposition may be carried out in bacteria, using a donor strain carrying the transposon either on a vector or as integrated copies in the genome. A target of interest is introduced into the transposer host where it receives integrations. Targets bearing insertions are then recovered from the host by genetic selection. Similarly, eukaryotic host cells, such as yeast, plant, insect, or mammalian cells, can be used to carry out the transposon mediated integration of an activation vector into a fragment of genomic DNA.

Isolation of mRNA and cDNA Produced from Activated Endogenous Genes

In additional embodiments, the present invention is directed to methods for isolating genes, particularly genes contained within the genome of a eukaryotic cell, that are activated using the vectors of the invention. These methods exploit the structure of the mRNA molecules produced using the non-targeted gene activation vectors of the invention. The methods of the invention described herein allow virtually any activated gene to be isolated, regardless of whether it has been previously isolated and characterized, and regardless of whether it has a known biological activity. This is made possible by the nature of the chimeric transcripts produced from the integrated vectors of the present invention. Using methods described herein, activation vectors can be integrated into the genome of a cell. Typically, the activation vectors, however, are integrated into the genome of many cells to produce a library of unique integration events. Each member of the library contains the vector located at a unique integration site(s), and potentially contains an activated endogenous gene. Gene activation occurs when the activation vector integrates upstream of the 3'-most exon of an endogenous gene and in an orientation capable of allowing transcription from the vector to proceed through the endogenous gene. The integration site may be in an intron or exon of the endogenous gene, or may be upstream of the transcription start site of the gene. Following integration, the activation constructs are designed to produce a transcript capable of splicing from an exon encoded by the activation vector to an exon encoded by the endogenous gene. As a result, a chimeric message is produced that contains the vector exon linked to the exons from an endogenous gene, wherein the endogenous exons are derived from the region located downstream of the vector integration site. The structure of this chimeric transcript can be exploited for gene discovery purposes. For example, the chimeric transcripts can be rapidly isolated to use as probes (to isolate the full length cDNA or genomic copy of the gene or to characterize the gene) or for direct sequencing and/or characterization.

To isolate the chimeric transcripts activated by vector insertion, cDNA is produced from a library member containing the activation event. It is also possible to isolate chimeric transcripts from pools of library members in order to increase the through-put of the procedure. cDNA can then be produced from the mRNA harvested from the activated cells. Alternatively, total RNA may be used to produce cDNA. In either case, first strand synthesis can be carried out using an oligo dT primer, an oligo dT/poly(A) signal primer, or a random primer. To facilitate cloning of the cDNA product, a poly dT based primer can be used with the structure: 5'-Primer X(dT)₁₋₁₀₀-3'. The oligo dT/poly(A) signal primer can have the structure 5'-(dT)₁₀₋₃₀-Primer X-N₀₋₆-TTTATT-3'. The random primer can have the structure: 5'-(Primer X)NNNNNN-3'. In each primer, Primer X is any sequence that can be used to subsequently PCR amplify target nucleic acid molecules. Where the activated gene amplification product is to be cloned, it is useful to include one or more restriction sites within the primer X sequence to facilitate subsequent cloning. Other primers recognized by those skilled in the art can be used to create first strand cDNA products, including primers that lack a Primer X region.

In accordance with the invention, the primers may be conjugated with one or more hapten molecules to facilitate subsequent isolation of nucleic acid molecules (e.g., first and/or second strand cDNA products) comprising such primers. After the primer becomes associated with the nucleic acid molecule (via incorporation during cDNA synthesis), selective isolation of the molecule containing the haptenylated primer may be accomplished using a corresponding ligand which specifically interacts with and binds to the hapten via ligand-hapten interactions. In preferred such aspects, the ligand may be bound to, for example, a solid support. Once bound to the solid support, the molecules of interest (haptenylated primer-containing nucleic acid molecules) can be separated from contaminating nucleic acids and other materials by washing the support matrix with a solution, preferably a buffer or water. Cleavage of one or more of the cleavage sites within the primer, or by treatment of the solid support containing

the nucleic acid molecule with a high ionic strength elution buffer, then allows for removal of the nucleic acid molecule of interest from the solid support.

Preferred solid supports for use in this aspect of the invention include, but are not limited to, nitrocellulose, diazocellulose, glass, polystyrene, polyvinylchloride, polypropylene, polyethylene, dextran, Sepharose, agar, starch, nylon, latex beads, magnetic beads, paramagnetic beads, superparamagnetic beads or microtitre plates and most preferably a magnetic bead, a paramagnetic bead or a superparamagnetic bead, that comprises one or more ligand molecules specifically recognizing and binding to the hapten molecule on the primer.

Particularly preferred hapten molecules for use on the primer molecules of the invention, include without limitation: (i) biotin; (ii) an antibody; (iii) an enzyme; (iv) lipopolysaccharide; (v) apotransferrin; (vi) ferrotansferrin; (vii) insulin; (viii) cytokines (growth factors, interleukins or colony-stimulating factors); (ix) gp120; (x) β -actin; (xi) LFA-1; (xii) Mac-1; (xiii) glycophorin; (xiv) laminin; (xv) collagen; (xvi) fibronectin; (xvii) vitronectin; (xviii) integrins $\alpha_v\beta_1$ and $\alpha_v\beta_3$; (xix) integrins $\alpha_3\beta_1$, $\alpha_4\beta_1$, $\alpha_4\beta_7$, $\alpha_5\beta_1$, $\alpha_v\beta_1$, $\alpha_{mb}\beta_3$, $\alpha_v\beta_3$ and $\alpha_v\beta_6$; (xx) integrins $\alpha_1\beta_1$, $\alpha_2\beta_1$, $\alpha_3\beta_1$ and $\alpha_v\beta_3$; (xxi) integrins $\alpha_1\beta_1$, $\alpha_2\beta_1$, $\alpha_3\beta_1$, $\alpha_6\beta_1$, $\alpha_7\beta_1$ and $\alpha_6\beta_3$; (xxii) ankyrin; (xxiii) C3bi, fibrinogen or Factor X; (xxiv) ICAM-1 or ICAM-2; (xxv) spectrin or fodrin; (xxvi) CD4; (xxvii) a cytokine (e.g., growth factor, interleukin or colony-stimulating factor) receptor; (xxviii) an insulin receptor; (xxix) a transferrin receptor; (xxx) Fe^{+++} ; (xxxi) polymyxin B or endotoxin-neutralizing protein (ENP); (xxxii) an enzyme-specific substrate; (xxxiii) protein A, protein G, a cell-surface Fc receptor or an antibody-specific antigen; and (xxxiv) avidin and streptavidin. Particularly preferred is biotin.

Particularly preferred ligand molecules according to this aspect of the invention, which correspond in order to the above-described hapten molecules, include without limitation: (i) avidin and streptavidin; (ii) protein A, protein G, a cell-surface Fc receptor or an antibody-specific antigen; (iii) an enzyme-specific substrate; (iv) polymyxin B or endotoxin-neutralizing protein (ENP); (v) Fe^{+++} ; (vi) a transferrin receptor; (vii) an insulin receptor; (viii) a cytokine (e.g., growth

factor, interleukin or colony-stimulating factor) receptor; (ix) CD4; (x) spectrin or fodrin; (xi) ICAM-1 or ICAM-2; (xii) C3bi, fibrinogen or Factor X; (xiii) ankyrin; (xiv) integrins $\alpha_1\beta_1$, $\alpha_2\beta_1$, $\alpha_3\beta_1$, $\alpha_6\beta_1$, $\alpha_7\beta_1$ and $\alpha_6\beta_3$; (xv) integrins $\alpha_1\beta_1$, $\alpha_2\beta_1$, $\alpha_3\beta_1$ and $\alpha_6\beta_3$; (xvi) integrins $\alpha_3\beta_1$, $\alpha_4\beta_1$, $\alpha_4\beta_7$, $\alpha_5\beta_1$, $\alpha_6\beta_1$, $\alpha_{mb}\beta_3$, $\alpha_6\beta_3$ and $\alpha_6\beta_6$; (xvii) integrins $\alpha_6\beta_1$ and $\alpha_6\beta_3$; (xviii) vitronectin; (xix) fibronectin; (xx) collagen; (xxi) laminin; (xxii) glycophorin; (xxiii) Mac-1; (xxiv) LFA-1; (xxv) β -actin; (xxvi) gp120; (xxvii) cytokines (growth factors, interleukins or colony-stimulating factors); (xxviii) insulin; (xxix) ferrotransferrin; (xxx) apotransferrin; (xxxi) lipopolysaccharide; (xxxii) an enzyme; (xxxiii) an antibody; and (xxxiv) biotin. Particularly preferred, for use with biotinylated primers of the invention, are avidin and streptavidin.

Following first strand synthesis, second strand cDNA synthesis may be carried out using a primer specific for the vector encoded exon. This creates double stranded cDNA from all transcripts that were derived from the vector encoded promoter. All cellular mRNA (and cDNA) produced from endogenous promoters remains single stranded since the transcript lacks a vector exon at its 5' end. Once second strand synthesis is carried out, the cDNA may be digested with a restriction enzyme, cloned into a vector, and propagated.

To facilitate cloning, cDNA molecules containing the vector exon are amplified by PCR using a primer specific for the vector exon and a primer specific for the first strand cDNA primer (e.g. Primer X). PCR amplification results in the production of variable length DNA fragments representing different locations of priming during first strand synthesis and/or amplification of multiple chimeric transcripts from different genes. These amplification products can be cloned into plasmids for characterization, or can be labeled and used as a probe.

Other amplification techniques, such as linear amplification using RNA polymerase (Van Gelder, *Proc. Natl. Acad. Sci. USA* 87:1663-1667 (1990); Eberwine, *Methods* 10:283-288 (1996)), can be used. For example, when linear amplification by RNA polymerase is used, a promoter (e.g. T7 promoter) can be placed on the vector exon. As a result, gene activated transcripts will contain the

promoter sequence at the 5' end of the transcript. Alternatively, a promoter can be ligated onto the cDNA molecule following first strand and second strand synthesis. Using either strategy, RNA polymerase is then incubated with cDNA in the presence of ribonucleotide triphosphates to create RNA transcripts from the cDNA. These transcripts are then reverse transcribed to produce cDNA. Since RNA polymerase can create several thousand transcripts from a single cDNA molecule, and since each of these transcripts can be reverse transcribed into cDNA, a large amplification can be achieved. As with PCR, amplification with RNA polymerase can facilitate cloning of activated genes. Other types of amplification strategies are also possible.

In another embodiment, the vector exon containing cDNA molecules are isolated without amplification. This may be useful in instances where biases occur during amplification (for example, when one DNA fragment amplifies more efficiently than another). To produce cDNA enriched for tagged messages, RNA is isolated from the activation library. A primer (e.g. a random hexamer, oligo(dT), or hybrid primers containing a primer linked to poly(dT) or a random nucleotides) is annealed to the RNA and used to direct first strand synthesis. The first strand cDNA molecules are then hybridized to a primer specific for the vector encoded exon. This primer directs second strand synthesis. Following second strand synthesis, the cDNA may be digested with restriction enzymes that cut in the vector exon and in the first strand primer (e.g. in Primer X - see above). The second strand products may then be cloned into a useful vector to allow them to be propagated.

It will be apparent to one of ordinary skill in view of the description contained herein that the cDNA products made according to the methods of the invention may also be cloned into a cloning vector suitable for transfection or transformation of a variety of prokaryotic (bacterial) or eukaryotic (yeast, plant or animal including human and other mammalian) cells. Such cloning vectors, which may be expression vectors, include but are not limited to chromosomal-, episomal- and virus-derived vectors, e.g., vectors derived from bacterial plasmids

or bacteriophages, and vectors derived from combinations thereof, such as cosmids and phagemids, BACs, MACs, YACs, and the like. Other vectors suitable for use in accordance with this aspect of the invention, and methods for insertion of DNA fragments therein and transformation of host cells with such cloning vectors, will be familiar to those of ordinary skill in the art.

Removal of Unspliced Transcription Products

In some instances, the activation vector will integrate into the genome in a region lacking genes. Alternatively, it may integrate into a region containing a gene(s), but be oriented in a manner that results in the transcription of the non-coding strand. In each of these instances, gene activated transcripts are produced that contain normally untranscribed DNA sequences next to the vector encoded exon. These sequences would complicate identification and analysis of novel genes. Therefore, it would be advantageous to selectively remove these genomic molecules.

To remove cDNA molecules that contain a vector encoded intron, the double strand cDNA is treated with a restriction enzyme that recognizes a sequence located in the vector encoded intron. Preferably, the restriction enzyme creates an overhang that is different from the overhang produced by cleavage of the vector exon. This ensures the cloning of only activated genes by preventing the cleavage products from ligating into the cloning vector.

Recovery of Exon I from activated endogenous genes

To recover exon I from activated genes, specialized vectors can be used to create non-targeted gene activation libraries. In its simplest form, this vector contains, from 5' to 3', a promoter, an unpaired splice donor site, and a second promoter. The downstream promoter is oriented in the same direction as the upstream promoter. Upon integration upstream of an endogenous gene, this type of vector produces two types of transcripts. The first transcript contains the vector exon joined to exon II of the endogenous gene. Methods for isolating this

transcript are described above. The second transcript contains the upstream region of the endogenous gene followed by exon I joined to exon II and other downstream exons from the endogenous gene (Figure 6).

Using a two step process, exon I can be recovered from cells containing the integrated vector. First, vector exon containing transcripts (i.e. Transcript type #1, Figure 13) are isolated using the methods described above. Once isolated, the 5' end of the transcript including exon II can be sequenced to determine the sequence of the flanking endogenous exons. Second, once the sequence of the flanking endogenous exons is known, PCR primers capable of annealing to exon II (or a downstream exon) of the activated gene can be developed. These primers can be used to amplify exon I from Transcript #2 (Figure 13) using a modified form of inverse PCR (Zeiner, M., *Biotechniques* 17(6):1051-1053 (1994)). Briefly, amplification of exon I from the endogenous gene is achieved by carrying out first strand cDNA synthesis with a gene specific primer, based on the sequence information determined above. Second strand synthesis can be carried out using *E. coli* DNA polymerase I under conditions well known to those skilled in the art. The double strand cDNA is then digested with a restriction enzyme that cleaves at least once in the endogenous gene upstream of the first strand cDNA primer, and that does not cleave in the vector exon. Following digestion, the cDNA is self ligated to produce circular molecules. Using inverted PCR primers that anneal in the endogenous gene upstream of the restriction/circularization site, amplification by PCR produces a DNA product containing exon I sequences from the endogenous gene.

Method for Selecting Cells Containing Higher Levels of Gene Activated Transcripts/Protein

In several embodiments of the disclosed invention, the activation vector contains an amplifiable marker (e.g. DHFR) and a viral origin of replication (e.g. EBV ori P). In other embodiments, an amplifiable marker and viral origin of replication are present on a cloning vector containing a cloned fragment of

genomic DNA. In yet another embodiment, the activation vector contains one element (e.g. DHFR) and a cloning vector carrying a genomic insert contains the other element (e.g. Ori P). Regardless of the initial location of the amplifiable marker and viral origin, the elements are combined on the same DNA molecule prior to or during introduction into a host cell.

In addition to the cis-acting elements, a trans-acting viral protein is generally required for efficient replication of the episomes. Examples of trans-acting viral proteins include EBNA-1 and SV40 T antigen. To promote efficient replication of episomes, the trans-acting viral protein can be expressed from the episome. Thus, the viral trans-acting protein may be expressed from the transposing activation vector, or may be positioned on the backbone of the cloning vector. Alternatively, the trans-acting viral protein may be expressed by the eukaryotic host cells into which the episome is introduced.

Once the amplifiable marker and viral origin of replication are on the same molecule and present in a host cell expressing the appropriate viral replication protein(s), the copy number of the episome can be increased. To increase the copy number of the episome, the cells can be placed under the appropriate selection. For example, if DHFR is present on the episome, methotrexate may be added to the culture. The selective agent may be applied at relatively high concentrations to isolate cells in the population that already have a high episome copy number. Alternatively, the selective agent may be applied at lower concentrations, and periodically increased in concentration. Two-fold increases in drug concentration will result in step-wise increases in copy number.

To reduce the frequency of non-specific drug resistance (i.e. drug resistance that is not associated with increased copy number of the episome), more than one amplifiable marker can be placed on the vector. Inclusion of multiple amplifiable markers on the episome allows cells to be selected with multiple drugs (either simultaneously or sequentially). Since non-specific drug resistance is a relatively rare event, the probability of a cell developing non-specific drug resistance to multiple drugs is exceedingly rare. Thus, the presence of multiple

amplifiable markers on the episome facilitates isolation of cells that have a high episome copy number.

Amplification of episome copy number increases the number of transcripts derived from the vector activated gene. This, in turn, facilitates isolation of cDNA molecules derived from the activated gene. Furthermore, amplification of episome copy number can dramatically increase protein expression from the activated gene. Higher levels of protein production facilitate generation of proteins for bioassay screening, cell assay screening, and manufacturing purposes.

As a result of the highly desirable characteristics described above, vectors containing a viral origin of replication and an amplifiable marker, and the use of these vectors to rapidly amplify the copy number of episomal vectors, represent a break through that extends beyond the scope of activating expression of genes present in genomic DNA. For example, these vectors can be used to over-express cDNA encoded genes to produce high levels of protein expression without the need to integrate the gene into a host cell genome with an amplifiable marker. Furthermore, like amplification of chromosomal sequences, cell possessing several hundred to several thousand episomal copies of the vector can be isolated and maintained in culture. Thus, the vectors described herein, and their uses, allow high levels of cloned genomic DNA to be propagated in mammalian cells, facilitate isolation of cDNA copies of genes present on the vector as genomic inserts, and maximize protein production from cloned cDNA and genomic copies of eukaryotic genes.

Other suitable modifications and adaptations to the methods and applications described herein will be readily apparent to one of ordinary skill in the relevant arts and may be made without departing from the scope of the invention or any embodiment thereof. Having now described the present invention in detail, the same will be more clearly understood by reference to the following examples, which are included herewith for purposes of illustration only and are not intended to be limiting of the invention

EXAMPLES

Example 1: Transfection of Cells for Activation of Endogenous Gene Expression

Method: Construction of pRIG-1

Human DHFR was amplified by PCR from cDNA produced from HT1080 cells by PCR using the primers DHFR-F1

(5' TCCTTCGAAGCTTGTCATGGTTGGTTCGCTAAACTGCAT 3') (SEQ ID NO:1) and DHFR-R1 (5' AAACCTTAAGATCGATTAATCATTC-TTCTCATATACTTCAA 3') (SEQ ID NO:2), and cloned into the T site in pTARGET™ (Promega) to create pTARGET:DHFR. The RSV promoter was isolated from PREP9 by digestion with *NheI* and *XbaI* and inserted into the *NheI* site of pTARGET:DHFR to create pTgT:RSV+DHFR. Oligonucleotides JH169 (5' ATCCACCATGGCTACAGGTGAGTACTCG 3') (SEQ ID NO:3) and JH170 (5' GATCCGAGTACTCACCTGTAGCCATGGTGGATTAA 3') (SEQ ID NO:4) were annealed and inserted into the I-Ppo-I and *NheI* sites of pTgT:RSV+DHFR to create pTgT:RSV+DHFR+Exl. A 279 bp region corresponding to nucleotides 230-508 of pBR322 was PCR amplified using primers Tet F1 (5' GGCGAGATCTAGCGCTATATGCGTTGATGCAAT 3') (SEQ ID NO:5) and Tet F2 (5' GGCCAGATCTGCTACCTTAAGAGAGCCG-AAACAAGCGCTCATGAGCCCGAA 3') (SEQ ID NO:6). Amplification products were digested with *BglII* and cloned into the *BamHI* site of pTgT:RSV+RSV+DHFR+Exl to create pRIG-1.

Transfection -- Creation of pRIG-1 Gene Activation Library in HT1080 Cells

To activate gene expression, a suitable activation construct is selected from the group of constructs described above. The selected activation construct is then introduced into cells by any transfection method known in the art.

Examples of transfection methods include electroporation, lipofection, calcium phosphate precipitation, DEAE dextran, and receptor mediated endocytosis. Following introduction into the cells, the DNA is allowed to integrate into the host cell's genome via non-homologous recombination. Integration can occur at spontaneous chromosome breaks or at artificially induced chromosomal breaks.

Method: Transfection of human cells with pRIG1. 2×10^9 HH1 cells, an HPRT⁻ subclone of HT1080 cells, was grown in 150 mm tissue culture plates to 90% confluency. Media was removed from the cells and saved as conditioned media (see below). Cells were removed from the plate by brief incubation with trypsin, added to media/10% fetal bovine serum to neutralize the trypsin, and pelleted at 1000 rpm in a Jouan centrifuge for 5 minutes. Cells were washed in 1X PBS, counted, and repelleted as above. The cell pellet was resuspended at 2.5×10^7 cells/ml final in 1X PBS (Gibco BRL Cat #14200-075). Cells were then exposed to 50 rads of γ irradiation from a ^{137}Cs source. pRIG1 (Fig. 14A-14B; SEQ ID NO:18) was linearized with *Bam*HI, purified with phenol/chloroform, precipitated with ethanol, and resuspended in PBS. Purified and linearized activation construct was added to the cell suspension to produce a final concentration of 40 $\mu\text{g/ml}$. The DNA/irradiated cell mixture was then mixed and 400 μl was placed into each 0.4 cm electroporation cuvettes (Biorad). The cuvettes were pulsed at 250 Volts, 600 μFarads , 50 Ohms using an electroporation apparatus (Biorad). Following the electric pulse, the cells were incubated at room temperature for 10 minutes, and then placed into $\alpha\text{MEM}/10\%\text{FBS}$ containing penicillin/streptomycin (Gibco/BRL). The cells were then plated at approximately 7×10^6 cells/150 mm plate containing 35 ml $\alpha\text{MEM}/10\%\text{FBS}/\text{penstrep}$ (33% conditioned media/67% fresh media). Following a 24 hour incubation at 37°C, G418 (Gibco/BRL) was added to each plate to a final concentration of 500 $\mu\text{g/ml}$ from a 60 mg/ml stock. After 4 days of selection, the media was replaced with fresh $\alpha\text{MEM}/10\%\text{FBS}/\text{penstrep}/500 \mu\text{g/ml}$ G418. The cells were then incubated for another 7-10 days and the culture supernatant assayed for the presence of new protein factors or stored at -80 °C for later

analysis. The drug resistant clones can be stored in liquid nitrogen for later analysis.

Example 2: Use of Ionizing Irradiation to Increase the Frequency and Randomness of DNA Integration

Method: HH1 cells were harvested at 90% confluency, washed in 1x PBS, and resuspended at a cell concentration of 7.5×10^6 cells/ml in 1X PBS. 15 μ g linearized DNA (pRIG-I) was added to the cells and mixed. 400 μ l was added to each electroporation cuvette and pulsed at 250 Volts, 600 μ Farads, 50 Ohms using an electroporation apparatus (Biorad). Following the electric pulse, the cells were incubated at room temperature for 10 minutes, and then placed into 2.5 ml α MEM/10%FBS/1X penstrep. 300 μ l of cells from each shock were irradiated at 0, 50, 500, and 5000 rads immediately prior to or at either 1 hour or 4 hours post transfection. Immediately following irradiation, the cells were plated onto tissue culture plates in complete medium. At 24 hours post plating, G418 was added to the culture to a final concentration of 500 μ g/ml. At 7 days post-selection, the culture medium was replaced with fresh complete medium containing 500 μ g/ml G418. At 10 days post selection, medium was removed from the plate, the colonies were stained with Coomassie Blue/90% methanol/10% acetic acid and colonies with greater than 50 cells were counted.

Example 3: Use of Restriction Enzymes to Generate Random, Semi-random, or Targeted Breaks in the Genome

Method: HH1 cells were harvested at 90% confluence, washed in 1x PBS, and resuspended at a cell concentration of 7.5×10^6 cells/ml in 1X PBS. To test the efficiency of integration, 15 μ g linearized DNA (PGK- β geo) was added to each 400 μ l aliquot of cells and mixed. To several aliquots of cells, restriction enzymes *Xba*I, *Not*I, *Hind*III, *Ipp*oI (10-500 units) were then added to separate cell/DNA mixture. 400 μ l was added to each electroporation cuvette and pulsed at 250 Volts, 600 μ Farads, 50 Ohms using an electroporation apparatus (BioRad).

Following the electric pulse, the cells were incubated at room temperature for 10 minutes, and then placed into 2.5 ml α MEM/10%FBS/IX penstrep. 300 μ l of 2.5 ml total cells from each shock were plated onto tissue culture plates in complete media. At 24 hours post plating, G418 was added to the culture to a final concentration of 600 μ g/ml. At 7 days post-selection, the media was replaced with fresh complete media containing 600 μ g/ml G418. At 10 days post selection, media was removed from the plate, the colonies were stained with Coomassie Blue/90% methanol/10% acetic acid and colonies with greater than 50 cells were counted.

Example 4: Amplification by Selecting for Two Amplifiable Markers Located on the Integrated Vector

Following integration of the vector into the genome of a host cell, the genetic locus may be amplified in copy number by simultaneous or sequential selection for one or more amplifiable markers located on the integrated vector. For example, a vector comprising two amplifiable markers may be integrated into the genome, and expression of a given gene (*i.e.*, a gene located at the site of vector integration) can be increased by selecting for both amplifiable markers located on the vector. This approach greatly facilitates the isolation of clones of cells that have amplified the correct locus (*i.e.*, the locus containing the integrated vector).

Once the vector has been integrated into the genome by nonhomologous recombination, individual clones of cells containing the vector integrated in a unique location may be isolated from other cells containing the vector integrated at other locations in the genome. Alternatively, mixed populations of cells may be selected for amplification.

Cells containing the integrated vector are then cultured in the presence of a first selective agent that is specific for the first amplifiable marker. This agent selects for cells that have amplified the amplifiable marker either on the vector or on the endogenous chromosome. These cells are then selected for amplification of the second selectable marker by culturing the cells in the presence of a second

selective agent that is specific for the second amplifiable marker. Cells that amplified the vector and flanking genomic DNA will survive this second selective step, whereas cells that amplified the endogenous first amplifiable marker or that developed non-specific resistance will not survive. Additional selections may be performed in similar fashion when vectors containing more than two (*e.g.*, three, four, five, or more) amplifiable markers are integrated into the cell genome, by sequential culturing of the cells in the presence of selective agents that are specific for the additional amplifiable markers contained on the integrated vector. Following selection, surviving cells are assayed for level of expression of a desired gene, and the cells expressing the highest levels are chosen for further amplification. Alternatively, pools of cells resistant to both (if two amplifiable markers are used) or all (if more than two amplifiable markers are used) of the selective agents may be further cultured without isolation of individual clones. These cells are then expanded and cultured in the presence of higher concentrations of the first selective agent (usually twofold higher). The process is repeated until the desired expression level is obtained.

Alternatively, cells containing the integrated vector may be selected simultaneously for both (if two are used) or all (if more than two are used) of the amplifiable markers. Simultaneous selection is accomplished by incorporating both selection agents (if two markers are used) or all of the selection agents (if more than two markers are used) into the selection medium in which the transfected cells are cultured. The majority of surviving cells will have amplified the integrated vector. These clones can then be screened individually to identify the cells with the highest expression level, or they can be carried as a pool. A higher concentration of each selective agent (usually twofold higher) is then applied to the cells. Surviving cells are then assayed for expression levels. This process is repeated until the desired expression levels are obtained.

By either selection strategy (*i.e.*, simultaneous or sequential selection), the initial concentration of selective agent is determined independently by titrating the agent from low concentrations with no cytotoxicity to high concentrations that

result in cell death in the majority of cells. In general, a concentration that gives rise to discrete colonies (e.g., several hundred colonies per 100,000 cells plated) is chosen as the initial concentration.

Example 5: Isolation of cDNAs Encoding Transmembrane Proteins

pRIG8R1-CD2 (Fig. 5A-5D; SEQ ID NO:7), pRIG8R2-CD2 (Fig. 6A-6C; SEQ ID NO:8), and pRIG8R3-CD2 (Fig. 7A-7C; SEQ ID NO:9) vectors contain the CMV immediate early gene promoter operably linked to an exon followed by an unpaired splice donor site. The exon on the vector encodes a signal peptide linked to the extra-cellular domain of CD2 (lacking an in frame stop codon). Each vector encodes CD2 in a different reading frame relative to the splice donor site.

To create a library of activated genes, 2×10^7 cells were irradiated with 50 rads from a ^{137}Cs source and electroporated with 15 μg of linearized pRIG8R1-CD2 (SEQ ID NO:7). Separately, this was repeated with pRIG8R2-CD2 (SEQ ID NO:8), and again with pRIG8R3-CD2 (SEQ ID NO:9). Following transfection, the three groups of cells were combined and plated into 150 mm dishes at 5×10^6 transfected cells per dish to create library #1. At 24 hours post transfection, library #1 was placed under 500 $\mu\text{g}/\text{ml}$ G418 selection for 14 days. Drug resistant clones containing the vector integrated into the host cell genome were combined, aliquoted, and frozen for analysis. Library #2 was created as described above, except that 3×10^7 cells, 3×10^7 cells and 1×10^7 cells were transfected with pRIG8R1-CD2, pRIG8R2-CD2, and pRIG8R3-CD2, respectively.

To isolate cells containing activated genes encoding integral membrane proteins, 3×10^6 cells from each library were cultured and treated as follows:

- Cells were trypsinized using 4 mls of Trypsin- EDTA.
- After the cells had released, the trypsin was neutralized by addition of 8 ml of alpha MEM/10% FBS.

- The cells were washed once with sterile PBS and collected by centrifugation at 800 x g for 7 minutes.
- The cell pellet was resuspended in 2ml of alpha MEM/10% FBS. 1 ml was used for sorting while the other 1 ml was replated in alpha MEM/10% FBS containing 500 µg/ml G-418, expanded and saved.
- The cells used for sorting were washed once with sterile alpha MEM/10% FBS and collected by centrifugation at 800 x g for 7 minutes.
- The supernatant was removed and the pellet resuspended in 1 ml of alpha MEM/10% FBS. 100 µl of these cells was removed for staining with the isotype control.
- 200 µl of Anti-CD2 FITC (Pharmingen catalog # 30054X) was added to the 900 µl of cells while 20 µl of the Mouse IgG₁ isotype control (Pharmingen catalog # 33814X) was added to the 100 µl of cells. The cells were incubated, on ice, for 20 minutes.
- To the tube that contained the cells stained with the Anti-Human CD2 FITC, 5 ml of PBS/1% FBS were added. To the isotope control, 900 µl of PBS/1% FBS were added. The cells were collected by centrifugation at 600 x g for 6 minutes.
- The supernatant from the tubes was removed. The cells that had been stained with the isotype control were resuspended in 500 µl of alpha MEM/10% FBS, and the cells that had been stained with anti-CD2- FITC were resuspended in 1.5 ml alpha MEM/10% FBS.

Cells were sorted through five sequential sorts on a FACS Vantage Flow Cytometer (Becton Dickinson Immunocytometry Systems; Mountain View, CA). In each sort, the indicated percentage of total cells, representing the most strongly fluorescent cells (see below) were collected, expanded, and resorted. HT1080

cells were sorted as a negative control. The following populations were sorted and collected in each sort:

	Library #1	Library #2	Library #3
Sort #1	500,000 cells collected (top 10%)	100,000 cells collected (top 10%)	40,000 cells collected (top 10%)
Sort #2	300,000 cells collected (top 5%)	220,000 cells collected (top 11%)	14,000 cells collected (top 5%)
Sort #3	90,000 cells collected (top 5%)	40,000 cells collected (top 10%)	120,000 cells collected (top 10%)
Sort #4	600,000 cells collected (top 40%)	(a) 6,000 cells collected (top 5%); (b) 10,000 cells collected (next 5%)	280,000 cells collected (top 13%)
Sort #5	(a) 260,000 cells collected (top 10%); (b) 530,000 cells collected (next 25%)	(a) from group (a) of sort #4, 100,000 cells collected (top 10%), and 350,000 cells collected (next 35%); (b) from group (b) of sort #4, 120,000 cells collected (top 10%)	(Not done)

Cells from each of the final sorts for each library were expanded and stored in liquid nitrogen.

Isolation of activated genes from FACS-sorted cells

Once cells had been sorted as described above, activated endogenous genes from the sorted cells were isolated by PCR-based cloning. One of ordinary skill will appreciate, however, that any art-known method of cloning of genes may be equivalently used to isolate activated genes from FACS-sorted cells.

Genes were isolated by the following protocol:

(1) Using PolyATract System 1000 mRNA isolation kit (Promega), mRNA was isolated from 3×10^7 CD2+ cells (sorted 5 rounds by FACS, as described above) from libraries #1 and #2.

(2) After mRNA isolation, the concentration of mRNA was determined by diluting 0.5 μ l of isolated mRNA into 99.5 μ l water and measuring OD₂₆₀. 21 μ g of mRNA were recovered from the CD2+ cells.

(3) First strand cDNA synthesis was then carried out as follows:

(a) While the PCR machine was holding at 4°C, first strand reaction mixtures were set up by sequential addition of the following components:

41 μ l DEPC-treated ddH₂O

4 μ l 10mM each dNTP

8 μ l 0.1 MDTT

16 μ l 5x MMLV first strand buffer (Gibco-BRL)

5 μ l (10pmol/ μ l) of the consensus polyadenylation site primer GD.R1 (SEQ ID NO:10)*

1 μ l RNAsin (Promega)

3 μ l (1.25 μ g/ μ l) mRNA.

*Note: GD.R1, 5'TTTTTTTTTTTTCGTCAGCGCCGCATCNNNNNTTT-ATT 3' (SEQ ID NO:10), is a "Gene Discovery" primer for first strand cDNA synthesis of mRNA; this primer is designed to anneal to the poly-adenylation signal AATAAA and downstream poly-A region. This primer will introduce a *NotI* site into the first strand.

Once samples had been made up, they were incubated as follows:

(b) 70° for 1 min.

(c) 42° hold.

2 μ l of 400 U/ μ l SuperScript II (Gibco-BRL, Rockville, MD) was then added to each sample, to give a final total volume of 82 μ l. After approximately three minutes, samples were incubated as follows:

- (d) 37° for 30 min.
- (e) 94° for 2 min.
- (f) 4° for 5 min.

2 μ l of 20 U/ μ l RNase-IT (Stratagene) was then added to each sample, and samples were incubated at 37° for 10 min.

- (4) Following first strand synthesis, cDNA was purified using a PCR cleanup kit (Qiagen) as follows:

- (a) 80 μ l of the first strand reaction were transferred to a 1.7 ml siliconized eppendorf tube and adding 400 μ l of PB.
- (b) Samples were then transferred to a PCR clean-up column and centrifuged for two minutes at 14,000 RPM.
- (c) Columns were then disassembled, flowthrough decanted, 750 of μ l PE were added to pellets, and tubes were centrifuged for two minutes at 14,000 RPM.
- (d) Columns were disassembled and flowthrough decanted, and tubes then centrifuged for two minutes at 14,000 RPM to dry resin.

- (e) cDNA was then eluted using 50 μ l of EB through transferring column to a new siliconized eppendorf tube which was then centrifuged for two minutes at 14,000 RPM.

(5) Second strand cDNA synthesis was then carried out as follows:

- (a) Second strand reaction mixtures were set up at RT, through the sequential addition of the following components:

ddH ₂ O	55 μ l
10 x PCR buffer	10 μ l
50 mM MgCl ₂	5 μ l
10 mM dNTPs	2 μ l
25 pmol/ μ l RIG.751-Bio*	4 μ l
25 pmol/ μ l GD.R2**	4 μ l
First strand product	20 μ l

*Note: RIG.F751-Bio, 5' Biotin-CAGATCACTAGAAGCTTTATTGCGG 3' (SEQ ID NO:11), anneals at the cap-site of the transcript expressed from pRIG vectors.

**Note: GD.R2, 5' TTTTCGTCAGCGGCCGCATC 3' (SEQ ID NO:12), is a primer used to PCR amplify cDNAs generated using primer GD.R1 (SEQ ID NO:10). GD.R2 is a sub-sequence of GD.R1 with matching sequence up to the degenerate bases preceding the polyA signal sequence.

(b) Start second strand synthesis:

94°C for 1 min;

add 1 µl *Taq* (5U/µl, Gibco-BRL);

add 1 µl Vent DNA pol (0.1U/µl, New England Biolabs).

(c) Incubate at 63°C for 2 min.

(d) Incubate at 72°C for 3 min.

(e) Repeat step (b) four times.

(f) Incubate at 72°C for 6 min.

(g) Incubate at 4°C (hold)

(h) END

(6) 200 µl of 1 mg/ml Streptavidin-Paramagnetic Particles (SA-PMP) were then prepared by washing three times with STE.

(7) The products of the second strand reaction were added directly to the SA-PMPs and incubated at RT for 30 minutes.

(8) After binding, SA-PMPs were collected through the use of the magnet, and flowthrough material recovered.

(9) Beads were washed three times with 500 µl STE.

(10) Beads were resuspended in 50 µl of STE and collected at the bottom of the tube using the magnet. STE supernatant was then carefully pipetted off.

(11) Beads were resuspended in 50 µl of ddH₂O and placed into a 100°C water bath for two minutes, to release purified cDNA from PMPs.

- (12) Purified cDNA was recovered by collecting PMPs on the magnet and carefully removing the supernatant containing the cDNA.
- (13) Purified products were transferred to a clean tube and centrifuged at 14,000 RPM for two minutes to remove all of the residual PMPs.
- (14) A PCR reaction was then carried out to specifically amplify RIG activated cDNAs, as follows:

- (a) PCR reaction mixtures were set up at RT, through the sequential addition of the following components:

H ₂ O	59 µl
10 x PCR buffer	10 µl
50 mM MgCl ₂	5 µl
10 mM dNTPs	2 µl
25 pmol/µl RIG.F781*	2 µl
25 pmol/µl GD.R2	2 µl
second strand product	20 µl

*Note: RIG.F781, 5' ACTCATAGGCCATAGAGGCCTATCACAG-TTAAATTGCTAACGCAG 3' (SEQ ID NO:13), anneals downstream of GD.F1 GD.F3, GD.F5-Bio, and RIG.F751-Bio, and adds an *Sfi*I site for 5' cloning of cDNAs. This primer is used in nested PCR amplification of RIG Exon1 specific second strand cDNAs.

- (b) Start thermal cycler:
94°C for 3 min;
add 1 µl of *Taq* (5U/µl; Gibco-BRL);
add 1 µl of 0.1U/µl Vent DNA polymerase (New England Biolabs)

PCR was then carried out by 10 cycles of steps (c) to (e):

- (c) 94°C for 30 sec.
- (d) 60°C for 40 sec.
- (e) 72°C for 3 min.

PCR was then completed by carrying out the following steps:

- (f) 94°C for 30 sec.
- (g) 60°C for 40 sec.
- (h) 72°C for 3 min.
- (i) 72°C + 20 sec each cycle for 10 cycles
- (j) 72°C for 5 min
- (k) 4°C hold.

(15) After elution of library material with 50 µl EB, samples were digested by adding 10 µl of NEB Buffer 2, 40 µl of dH₂O and 2 µl of *Sfi*I and digesting for 1 hour at 50°C, to cut the 5' end of the cDNA at the *Sfi*I site encoded by the forward primer (RIG.F781; SEQ ID NO:13).

(16) Following *Sfi*I digestion, 5 µl of 1M NaCl and 2 µl of *Not*I were added to each sample, and samples digested for one hour at 37°C, to cut the 3' end of the cDNA at the *Not*I site encoded by the first strand primer (GD.R1; SEQ ID NO:10).

(17) The digested cDNA was then separated on a 1% low melt agarose gel. cDNAs ranging in size from 1.2Kb to 8Kb were excised from the gel.

(18) cDNA was recovered from the excised agarose gel using Qiaex II Gel Extraction (Qiagen). 2 µl of cDNA (approximately 30mg) was ligated to 7µl (35ng) of pBS-HSB (linearized with *Sfi*I/*Not*I) in a total volume of

5

- 10

# of 96-Well Plates:	1 Plate	2 Plates	3 Plates	4 Plates
Total # of 12.5 μ l PCR rxns:	96	192	288	384
dH_2O	755 μ l	1.47 ml	2.20 ml	2.94 ml
5X PCR Premix-4	250 μ l	500 μ l	750 μ l	1.0 ml
F Primers premix (25 pmol/ μ l)	10 μ l	20 μ l	30 μ l	40 μ l
R Primers premix (25 pmol/ μ l)	10 μ l	20 μ l	30 μ l	40 μ l
RNase-In Cocktail	3.2 μ l	6.3 μ l	9.6 μ l	12.8 μ l
Taq Polymerase (5 U/ μ l)	3.2 μ l	6.3 μ l	9.6 μ l	12.8 μ l
Total Volume (ml)	1.01	2.02	3.03	4.04

- (d) 10 μ l of the master mix were dispensed into each well of the PCR reaction plate.
- (e) 2.5 μ l from each 100 μ l *E. coli* culture were transferred into the corresponding wells of the PCR reaction plate.
- (f) PCR was performed, using typical PCR cycle conditions of:
 - (i) 94°C/2min. (Bacterial lysis and plasmid denaturation)
 - (ii) 30 cycles of 92°C denaturation for 15 sec; 60°C primer annealing for 20 sec; and 72°C primer extension for 40 sec
 - (iii) 72°C final extension for 5 min.
 - (iv) 4°C hold.
- (g) Bromophenol blue was then added to the PCR reaction; samples were mixed, centrifuged, and then the entire reaction mix was loaded onto an agarose gel.

- 23) Of 200 clones screened, 78% were positive for the vector exon. 96 of these clones were grown as minipreps and purified using a Qiagen 96-well turbo-prep following the Qiagen Miniprep Handbook (April 1997).
- 24) Many duplicate clones were eliminated though simultaneous digestion of 2 µl of DNA with *NofI*, *Bam* HI, *XhoI*, *XbaI*, *HindIII*, *EcoRI* in NEB Buffer 3, in a total volume of 22 µl, followed by electrophoresis on a 1% agarose gel.

Results:

Two different cDNA libraries were screened using this protocol. In the first library (TMT#1), eight of the isolated activated genes were sequenced. Of these eight genes, four genes encoded known integral membrane proteins and six were novel genes. In the second library (TMT#2), 11 isolated activated genes were sequenced. Of these 11 genes, one gene encoded a known integral membrane protein, one gene encoded a partially sequenced gene homologous to an integral membrane protein, and nine were novel genes. In all cases where the isolated gene correspond to a characterized known gene, that gene was an integral membrane protein.

Exemplary significant alignments (obtained from GenBank) for genes isolated from each library are shown below:

TMT#1 Significant Alignments:

179761|gb|M76559|HUMCACNLB Human neuronal DHP-sensitive voltage-dependent, calcium channel alpha-2b subunit mRNA complete CDs.
Length = 3600

>gi|3183974|emb|Y10183|HSMEMD H.sapiens mRNA for MEMD protein
Length = 4235

TMT#2 Significant Alignments:

>gi|476590|gb|U06715|HSU06715 Human cytochrome B561, HCYTO B561, mRNA, partial CDs.
Length = 2463

>gi|2184843|gb|AA459959|AA459959 zx66c01.s1 Soares total fetus
Nb2HF8 9w Homo sapiens cDNA clone 796414 3' similar to
gb:J03171 INTERFERON-ALPHA RECEPTOR PRECURSOR (HUMAN);
Length = 431

Example 6: Activation of Endogenous Genes using a Poly(A) Trap Vector

HT1080 cells (1×10^7 cells) were irradiated with 50 rads using a ^{137}Cs source and electroporated with 15 μg linearized pRIG14 (Figure 29A-29B). Following transfection, the cells were plated into a 150 mm dish at 5×10^6 cells/dish. At 24 hours, puromycin was added to 3 $\mu\text{g}/\text{ml}$. The cells were incubated at 37°C for 12 days in the presence of 3 $\mu\text{g}/\text{ml}$ puromycin. The media was replaced every 5 days. At 12 days, the number of colonies was counted, and the cells were trypsinized and replated onto a new dish. The cells were grown to 90% confluency and harvested for frozen storage and gene isolation. Typically, 1000-3000 colonies were produced per 1×10^7 cells transfected.

Example 7: Activation of Endogenous Genes Using a Dual Poly(A) Trap/SAT Vector

1×10^7 HH1 cells (HPRT-minus HT1080 cells) were irradiated with 50 rads using a ^{137}Cs source and electroporated with 15 μg linearized pRIG-22. Following transfection, the cells were plated into a 150 mm dish at 5×10^6 cells/dish. At 24 hours, neomycin was added to 500 $\mu\text{g}/\text{ml}$ G481. The cells were incubated at 37°C for 4 days in the presence of 500 $\mu\text{g}/\text{ml}$ G418. The media was replaced with fresh media containing 500 $\mu\text{g}/\text{ml}$ G418 and AgThg and grown in the presence of both drugs for an additional 7 days. Alternatively, as a control for HPRT activity, the media was replaced with fresh media containing 500 $\mu\text{g}/\text{ml}$ G418 and HAT (available from Life Technologies, Inc., Rockville, MD, and used at manufacturer's recommended concentration) and grown in the presence of both drugs for an additional 7 days. At 12 days post transfection, the number of colonies was counted, and the cells were trypsinized and replated onto a new dish.

The cells were grown to 90% confluency and harvested for frozen storage and gene isolation. Typically, cells subjected to G418/AgThg selection produced 1000-3000 colonies per 1×10^7 cells transfected. In contrast, cells subjected to G418/HAT selection produced approximated 100 colonies per 1×10^7 cells transfected.

Example 8: Isolation of activated genes

Non-targeted gene activation vectors are integrated into the genome of a eukaryotic cells using the methods of the invention. By integrating the vector into multiple cells, a library is created in which cells are expressing different vector activated genes. RNA is isolated from these cells using a commercial RNA isolation kit. In this example, RNA is isolated from cells using Poly(A) Tract 1000 (Promega). The RNA is converted into cDNA, amplified, size fractionated, and cloned into a plasmid for analysis and sequencing. A brief description of this process is presented.

- 1) Place 4 ml GTC Extraction buffer (Poly(A) tract 1000 Kit- Promega) in a 15 ml polycarbonate screw cap tube and add 168 μ l 2-mercaptoethanol and place in a 70°C water bath.
- 2) Place 8 ml dilution buffer in a 15 ml polycarbonate screw cap tube for every pellet processed and add 168 μ l 2-mercaptoethanol and place in a 70°C water bath.
- 3) Remove from -80°C storage cell pellets (1×10^7 - 1×10^8 cells) containing non-targeted gene activation vector integrated into their genome. Pipette 4ml GTC Extraction buffer immediately onto cell pellet. Pipette up-and-down several times until the pellet is resuspended and transfer into a 15 ml snap cap polypropylene tube.
- 4) Add the 8 ml dilution buffer and mix by inversion.
- 5) Add 10 μ l (500 pmol) of the biotinylated oligo dT primer and mix.
- 6) Let sit at 70°C for 5 minutes inverting every couple of minutes to ensure even heating.

7) Centrifuge in a Sorvall HB-6 rotor at 7800 rpm (10k x g) at 25°C for 10 minutes. During this period of time wash 6 ml Streptavidin-Paramagnetic particles (SA-PMPs) 3x with 6 ml 0.5x SSC through use of the Poly(A) Tract system 1000 magnet.

8) After 3 washes resuspend the SA-PMPs in 6 ml 0.5 x SSC.

9) Pipette to remove the supernatant from the RNA prep and add to the resuspended SA-PMPs (Be careful when removing supernatant so that you do not disrupt the pellet).

10) Let the SA-PMP/RNA mix and incubate for 2 minutes at room temperature.

11) Capture the magnetic beads through use of the Poly(A) Tract system 1000 magnet. Note that it takes some time for all of the beads to pellet due to the high viscosity of the liquid.

12) Pour off the supernatant and resuspend the beads in 1.7 ml of 0.5 x SSC using a 2 ml pipette and transfer to a 2 ml screw cap tube.

13) Capture the SA-PMPs using the magnet and remove the supernatant by pipetting with a P1000.

14) Add 1.7 ml 0.5x SSC and invert the tube several times to mix.

15) Repeat steps 14 and 15 two more times.

16) Resuspend the SA-PMPs in 1 ml of nuclease free water and invert several times to mix.

17) Capture the SA-PMPs and pipette off the mRNA.

18) Place 0.5 ml of the mRNA into each of two siliconized eppendorf tubes and add 50 µl of DEPC-treated 3M NaOAc solution and 0.55 ml of isopropanol. Invert several times to mix and place at -20°C for at least 4 hours.

19) Centrifuge the mRNA for 10 minutes at max RPM (14 k).

20) Carefully pipette off the supernatants and wash pellets with 200 µl 80% ethanol through re-centrifugation for 2 minutes at 14K RPM. Note that the pellets are often brown or tan in color. This color results from residual SA-PMPs.

21) Remove wash and let pellets air dry for not more than 10 minutes at room temperature.

- 22) Resuspend pellets in 5 μ l each and combine into a single tube.
- 23) Centrifuge at 14K RPM for 2 minutes to remove the residual SA-PMPs and carefully remove the mRNA.
- 24) Determine the concentration of mRNA by diluting 0.5 μ l into 99.5 μ l water and measuring OD 260. Note that 1 OD 260 = 40 μ g RNA.
- 25) Set up first strand reaction for both the test sample and the negative control (HT1080) through the sequential addition of the following components while the PCR machine is holding at 4°C:

Step 1:

42 μ l DEPC-treated ddH₂O
4 μ l 10mM each dNTP
8 μ l 0.1 M DTT
16 μ l 5x MMLV 1st strand buffer
5 μ l (10pmol/ μ l) GDR1
1 μ l RNasin (Promega)
4 μ l (1.25 μ g/ μ l) mRNA.

Step 2: 70°/1 min

Step 3: 42°/hold

Step 4: After 1 minute add 2 μ l SUPERScript II® (Life Technologies, Inc.; Rockville, MD) and incubate at 37°C for 30 min

Step 5: 94°/2 min

Step 6: 4°/∞

Step 7: Add 2 μ l RNase and incubate at 37°C for 10 min

Step 8: 4°/∞

26) Analyze 8 μ l of cDNA on a 1% agarose gel to check for cDNA synthesis and purify remaining cDNA using the PCR cleanup kit from Qiagen by transferring the 70 μ l first strand reaction to a 1.5 ml siliconized eppendorf tube and adding 400 μ l PB.

27) Transfer to a PCR clean-up column and centrifuge 2 minutes at max RPM.

28) Disassemble column and pour out Flow through. Add 750 μ l PE and centrifuge 2 minutes at max RPM.

29) Disassemble column and pour out Flow through then centrifuge 2 minutes at max RPM to dry resin.

30) Elute using 50 μ l of EB through transferring column to a new siliconized eppendorf tube and centrifuging for 2 minutes at max RPM.

31) Second Strand cDNA synthesis set up at RT:

H ₂ O	8.5 μ l
10 X PCR buffer	5 μ l
50 mM MgCl ₂	2.5 μ l
10 mM dNTPs	1 μ l
25 pmol/ μ l GDF5Bio	10 μ l
25 pmol/ μ l GDR2	10 μ l
First strand product	15 μ l

Step 9: 94°C/1 min.

Step 10: 60°C/10 min.

Add 0.25 μ l *Taq* polymerase

Step 11: 60°C/2 min.

Step 12: 72°C/10 min.

Step 13: 94°C/1 min.

Step 14: min go to "Step 11" four more times

Step 15: 60°C/2 min

Step 16: 72°C/10 min

Step 17: END

32) Prepare 100 μ l of SA-PMPs by washing 3 x with STE and collection using a magnet. After the final wash, resuspend the beads in 150 μ l STE.

33) Purify the products of the second strand reaction using the PCR cleanup kit from Qiagen. Elute in 50 μ l EB and add the products of the second strand reaction to 150 μ l of the PMPs.

34) Mix gently at RT for 30 minutes.

35) After binding collect SA-PMPs through use of a magnet and recover flow through material (SAVE THIS MATERIAL!)

36) Wash the beads 3 x with 500 μ l STE and 1x with NEB 2 (1x).

37) Resuspend the beads in 100 μ l NEB 2 (1x).

38) Add 2 μ l *Sfi*I and digest at 50°C for 30 minutes with gentle mixing every 10 minutes.

39) Recover purified cDNA through use of a magnet and carefully removing the supernatant.

5 40) Transfer the products to a new tube and centrifuge at maximum RPM for 2 minutes to remove all of the beads.

41) Set up a PCR reaction to specifically amplify RAGE activated cDNAs:

H ₂ O	37 μ l
10 X PCR buffer	10 μ l
10 mM dNTPs	2 μ l
25 pmol/ μ l GDF 781	10 μ l
25 pmol/ μ l GDR2	10 μ l
Second strand product	25 μ l

Step 1: 94°C/2 min.

Step 2: 94°C/45 sec.

Step 3: 60°C/10 min.

Add 0.5 μ l *Taq* Polymerase

Step 4: 72°C/10 min.

Step 6: 60°C/2 min.

Step 7: 72°C/10 min.

Step 8: Cycle to step 5, 8 more times

Step 9: 94°C/45 sec.

Step 10: 60°C/2 min.

Step 11: 72°C/ 10 min. + 20 sec each cycle

Step 12: Cycle to step 9, 14 more times

Step 13: 72°C/ 5 min.

Step 14: 4°C hold

42) Check specificity of PCR amplification of HT1080 versus library material through analysis on a 1% agarose gel. If there is a high specificity of cDNA amplification, then use Qiagen PCR clean up kit to purify PCR products.

43) After elution of library material with 50 μ l EB add 10 μ l NEB2, 40 μ l dH₂O and 2 μ l *Sfi*I and digest for 1 hour at 50°C.

44) Add 5 μ l of 1 M NaCl and 2 μ l of *Not*I and digest for 1 hour at 37°C.

45) Prepare and run a 1% L.M. agarose gel and run library material on gel. After visualization of material, cut out fragments ranging in size from 500bp to 10 Kb.

46) Recover the library DNA from agarose using Qiaex II Gel Extraction Protocol (Qiagen) and elute DNA in 10 μ l EB. Ligate 5 μ l of this material to 4 μ l pBS-HSB (*Sfi*I/*Not*I) or pBS-SNS in a total volume of 10 μ l.

47) Transform *E. coli* with 0.5 μ l ligated DNA per 40 μ l cells.

48) Pick colonies, grow overnight in LB, isolate plasmids.

49) Analyze gene activated cDNA inserts by restriction digest and DNA sequencing.

Example 9: Isolation of Activated Genes from Subtracted cDNA Pools

Purified mRNAs from non-transfected HT1080 cells was prepared using the Poly-A Tract 1000 system (Promega), as described in Example 8 steps 1-24, and were biotinylated using EZ-Link™ Biotin LC-ASA reagent (Pierce), as follows:

1.) 25 μ l DEPC-treated dH₂O and 15 μ l containing 10 μ g of HT1080 mRNA was added into a siliconized microfuge tube and held on ice.

2.) Working under subdued light, 40 μ l of prepared LC-ASA stock reagent (1 mg/ml in 100% ethanol) was added into the reaction tube.

3.) A UV light (365 nm wavelength) was positioned 5 cm above the microfuge tube and used to irradiate the reaction mix for 15 minutes.

4.) Unlinked biotin reagent was removed from the labeled HT1080 mRNA by passing the reaction mix through an RNase-free MicroSpin P-30 column (BioRad), as prescribed by the manufacturer.

5 HT1080 cells were transfected with a poly(A) trap pRIG activation vector and grown under selective media to produce a population of drug resistant colonies, as described in Example 1. Purified mRNAs were prepared from the pooled colonies using the Promega Poly-A Tract 1000 system, as described in Example 8. First strand cDNA was prepared from 5 µg of this mRNA using oligo GD.R1(TTTTTTTTTTTCGTCAGCGCCGCATC>NNNTTTATT)(SEQ ID NO:10), as described in Example 8, Step 25. The reaction mix was passed through a Qiagen PCR Quick Clean-up column and the purified 1st strand cDNA was recovered in 100 µl EB.

10 The subtractive hybridization of biotinylated HT1080 mRNAs (subtractor population) and 1st strand cDNAs prepared from the superpool of pRIG-transfected colonies (target population) was performed as follows:

- 1.) 9 µg of biotinylated mRNA was added into a 0.5 ml microfuge tube containing 0.5 µg 1st strand cDNA.
- 2.) 1/100x volume of 10 mg/ml glycogen, 1/10x volume of 3 M sodium acetate, pH 5.5, and 2.6x volume of 100% ethanol were added into the tube and mixed.
- 20 3.) The tube was placed at -80°C for 1 hr, then spun in a refrigerated microfuge for 20 minutes.
- 4.) The pellet of precipitated nucleic acids was drained, washed once with 70% ethanol, then air-dried.
- 5.) The pellet was solvated in 5 µl HBS (50 mM HEPES, pH 7.6; 2 mM EDTA; 0.2% SDS; 500 mM NaCl) and overlaid with 5 µl light mineral oil, then heated to 95°C for 2 minutes followed by 68°C for 24 hours.
- 25 6.) The reaction mix was diluted with 100 µl HB (HBS without SDS) and extracted once with 100 µl chloroform to remove the oil.
- 7.) The diluted hybridization mix was added to 300 µl streptavidin-coated paramagnetic particles (Promega) which had been pre-washed 3x in 300 µl HB.
- 30

8.) The mix was incubated 10 minutes at room temperature and the SA-PMP's and bound Biotin-mRNA:DNA hybrids were removed from solution by magnetic capture.

9.) Steps 7 and 8 were repeated once.

10.) The cleared solution was subjected to one additional round of subtractive hybridization and magnetic removal of captured hybrids (Steps 1-9), with the following exceptions:

Step 6: the hybridization reaction was diluted with 2x PCR Buffer (40 mM Tris-HCl, pH 8.4; 100 mM KCl).

Step 7: PMPs were pre-washed in 1X PCR Buffer

The twice-subtracted 1st strand cDNA was used to generate 2nd strand cDNA by combining 45 μ l of 1st strand cDNA with 7 μ l dH₂O, 5 μ l 50 mM MgCl₂, 2 μ l premix of 10 mM each dNTP, 1 μ l 10x PCR Buffer, 20 μ l of 12.5 pmol/ μ l GD19F1-Bio (5' Biotin-CTCGTTTAGTGC GGCCGCTCAG-ATCACTGAATTCTGACGACCT) (SEQ ID NO:14), 20 μ l of 12.5 pmol/ μ l GD.R2 (TTTTCGTCAGCGGCCGCATC) (SEQ ID NO:12), and 0.5 μ l Taq Polymerase, with thermocycling as described in Example 8, Step 31. The second strand cDNA product was amplified and further processed for the production of an *E. coli*-based cDNA library, as described in Example 8, steps 32-49.

Example 10: Selective Capture of RIG-activated Transcripts

HT1080 cells were transfected with pRIG19 activation vector (Figure 30A-30C) and cultured for 2 weeks in selective media, as described in Example 6. Total RNA was prepared from a pellet comprised of 10⁸ cells using TRIzol® Reagent (Life Technologies, Inc.; Rockville, MD) following the manufacturer's protocol, and was dissolved in 720 μ l of DEPC-treated dH₂O (dH₂O^{DEPC}). Contaminating genomic DNA was eliminated from the RNA preparation by mixing 80 μ l NEB 10x Buffer 2, 8 μ l Promega RNasin, and 20 μ l

RQ1 Promega RNase-free DNase, incubating at 37°C for 30 minutes, extracting sequentially with equal volumes of phenol:chloroform (1:1) and chloroform, mixing with 1/10x volume sodium acetate (pH 5.5), precipitating the RNA with 2x volume of 100% ethanol, and solvating the dried RNA pellet in dH₂O^{DEPC} to a final concentration of 4.8 µg/µl.

mRNA transcripts derived from pRIG19-activated genes were selectively captured from the pool of total cellular RNAs by mixing in a 2 ml RNase-free microfuge tube 150 µl total RNA, 150 µl HBDEPC (50 mM HEPES, pH 7.6; 2 mM EDTA; 500 mM NaCl), 3 µl Promega RNasin, and 2.5 µl (25 pmol/µl) oligo GD19.R1-Bio (see Table 1), then incubating at 70°C for 5 minutes followed by 50°C for 15 minutes. One ml of Promega streptavidin coated paramagnetic particles (SA-PMPs) was magnetically captured and washed 3x each with 1.5 ml of 0.5x SSC, and the SA-PMPs were left without being resuspended. The warm oligo:RNA hybridization reaction was added directly into the tube containing the semi-dry SA-PMPs. After incubating for 10 minutes at room temperature the SA-PMPs were washed 3x with 1ml 0.5x SSC.

Table 1: Primer and Oligonucleotide Sequences

	Primer/Oligo Name	Sequence	SEQ ID NO:
Forward PCR Primers	GD19.F1-Bio	5' Biotin-CTCGTTTAGTGCGG-CCGCTCAGATCACTGAATTC TGACGACCT	14
	GD19.F2-Bio	5' Biotin-CTCGTTTAGTGCGG-CGCCAGATCACTGAATTCTG ACGACCT	15
	GD19.F2	GACCTACTGATTAAACGGCC-ATA	16

Reverse PCR Primers	GD.R1	TTTTTTTTTTTCGTCAGCG- GCCGCATCNNNNTTTATT	10
	GD.R2	TTTTCGTCAGCGGCCGCATC	12
mRNA Capture Oligo	GD19.R1-Bio	TCGTCAGAATTCAGTGAT- CT-3' Biotin	17

After the final magnetic capture, the SA-PMP's were suspended in 190 μ l dH₂ODEPCand incubated at 68°C for 15 minutes. PMPs were immobilized by exposure to a magnetic and the cleared solution containing RIG-activated transcripts was transferred to a microfuge tube. 63 μ l of captured RIG-activated transcript were transferred to a PCR tube where first and second strand cDNA synthesis was performed using PCR program "1+2CDNA", as follows:

- Step 1:* 4°C/ ∞ : Add into the PCR tube containing the RIG-activated transcripts 20 μ l 5x GibcoBRL RT Buffer, 1 μ l Promega RNasin, 10 μ l 100 mM DTT, 5 μ l dNTP premix at 10 mM each, 1 μ l oligo GD.R1 (see Table 1) at 25 pmol/ μ l.
- Step 2:* 70°C/3 minutes
- Step 3:* 42°C/10 minutes
- Step 4:* Add 2.5 μ l SUPERSCRIPT II® (Life Technologies, Inc.), then incubate at 37°C/1 hour
- Step 5:* 94°C/2 minutes
- Step 6:* 4°C/ ∞ .

To the 1st strand cDNA mix, 2 μ l of Stratagene RNase-It was added and the mixture was incubated at 37°C for 15 minutes. 600 μ l of Qiagen PB reagent was added to the reaction, then transferred to a Qiagen PCR clean-up column and processed according to the manufacturer's protocol. cDNA was eluted from the column in 50 μ l EB and transferred to a PCR tube. The second strand cDNA

reaction was performed using oligos GD19.F2-Bio (Table 1) and GD.R2 (Table 1) as described in Example 9. The second strand product was captured on Promega SA-PMPs as described in Example 9, with the exception that the final suspension of SA-PMPs was in 1x NEB 4 Buffer and the captured cDNAs were cleaved from the particles using restriction endonuclease Asc I. Amplification of the second strand cDNA products using oligos GD19.F2 and GD.R2, digestion of the amplified cDNAs using endonucleases *Sfi*I and *Not*I, and size selection of cDNAs prior to cloning were all performed as described in Example 9. The final cDNA cleanup was achieved by eluting the cDNA pool off a Qiagen PCR Cleanup column in 30 μ l EB. 11 μ l of cDNA was mixed with 4 μ l 5x GibcoBRL Ligase Buffer, 4 μ l pGD5 vector DNA previously prepared by digestion with *Sfi*I, *Not*I, and CIP. 1 μ l T4 DNA Ligase was added, and the reaction mix was incubated at 16°C overnight. 1 μ l of ligation reaction was used to transform electro-competent *E. coli* DH10B cells, which were subsequently plated on LB agar plates containing 12.5 μ g/ml chloramphenicol. Typically, 60 to 80 bacterial colonies were recovered per μ l of ligation mix transformed.

Example 11: Selective Capture of RIG-activated Transcripts

HT1080 cells were transfected with pRIG19 activation vector and cultured for 2 weeks in selective media, as described in Example 6. Total RNA was prepared from a pellet comprised of 10^8 cells using TRIzol® Reagent (Life Technologies, Inc.) following the manufacturer's protocol, and was dissolved in 720 μ l of DEPC treated dH₂O (dH₂O^{DEPC}). Contaminating genomic DNA was eliminated from the RNA preparation by mixing 80 μ l NEB 10x Buffer 2, 8 μ l Promega RNasin, and 20 μ l RQ1 Promega RNase-free DNase, incubating at 37°C for 30 minutes, extracting sequentially with equal volumes of phenol:chloroform (1:1) and chloroform, mixing with 1/10x volume sodium acetate (pH 5.5), precipitating the RNA with 2x volume of 100% ethanol, and solvating the dried RNA pellet in dH₂O/DEPC to a final concentration of 4.8 μ g/ μ l.

mRNA transcripts derived from pRIG19-activated genes were selectively captured from the pool of total cellular RNAs by mixing in a 2 ml RNase-free microfuge tube 150 μ l total RNA, 150 μ l HBDEPC (50 mM HEPES, pH 7.6; 2 mM EDTA; 500 mM NaCl), 3 μ l Promega RNasin, and 2.5 μ l (25 pmol/ μ l) oligo GD19.R1-Bio (see Table 1), then incubating at 70°C for 5 minutes followed by 50°C for 15 minutes. One ml of Promega streptavidin coated paramagnetic particles (SA-PMPs) was magnetically captured and washed 3x each with 1.5 ml of 0.5x SSC, and the SA-PMPs were left without being resuspended. The warm oligo:RNA hybridization reaction was added directly into the tube containing the semi-dry SA-PMPs. After incubating for 10 minutes at room temperature the SA-PMPs were washed 3x with 1ml 0.5x SSC. After the final magnetic capture the SA-PMP's were suspended in 190 μ l dH₂O^{DEPC} and incubated at 68°C for 15 minutes. PMPs were immobilized by exposure to a magnetic and the cleared solution containing RIG-activated transcripts was transferred to a microfuge tube. 63 μ l of captured RIG-activated transcript were transferred to a PCR tube where first and second strand cDNA synthesis was performed using PCR program "1+2CDNA", as follows:

Step 1: 4°C/ ∞ : Add into the PCR tube containing the RIG-activated transcripts 20 μ l 5x GibcoBRL RT Buffer, 1 μ l Promega RNasin, 10 μ l 100 mM DTT, 5 μ l dNTP premix at 10 mM each, 1 μ l oligo GD.R1 (see Table 1) at 25pmol/ μ l.

Step 2: 70°C/3 minutes

Step 3: 42°C/10 minutes

Step 4: Add 2.5 μ l SUPERSCRIPT II® (Life Technologies, Inc.), then incubate at 37°C/1 hour

Step 5: 94°C/2 minutes

Step 6: 60°C/ ∞ ; while holding temperature, the following were added: 2 μ l 50 mM MgCl₂, 1 μ l oligo GD19.F1-Bio (Table 1) at 25 pmol/ μ l, and 2 μ l Stratagene RNase-It.

After 10 minutes, 0.5 μ l *Taq* DNA Polymerase (Life Technologies, Inc.) was added and the cycling was continued:

Step 7: 72°C/10 minutes

Step 8: 4°C/∞.

The 100 μ l volume cDNA reaction mix was transferred to a 1.5 ml siliconized microfuge tube and extracted sequentially with equal volumes of phenol:chloroform (1:1) and chloroform, and the aqueous phase was transferred to a new tube and place in speed-vac for 5 minutes at 37°C. Restriction digestion of the cDNA was performed by adding 74 μ l dH₂O, 20 μ l NEB 10x Buffer 2, 2 μ l 1 mg/ml BSA, 4 μ l *Sfi*I and incubating at 50°C for 1 hour, then adding 10 μ l 1 M NaCl, 4 μ l *Not*I and incubating an additional 37°C for 1 hour. The reaction mix was extracted sequentially with equal volumes of phenol:chloroform (1:1) and chloroform, then cDNAs were precipitated by adding 1/100x volume 10 mg/ml glycogen, 1/30x volume 3 M sodium acetate (pH 7.5), 2x volume 100% absolute ethanol, and freezing at -80°C for 1 hour. The cDNA pellet was washed once with 70% ethanol and air dried for 15 minutes, then solvated in 5 μ l dH₂O, 1 μ l 10X NEB Ligase Buffer, 4 μ l pGD5 vector DNA previously prepared by digestion with *Sfi*I, *Not*I, and CIP. 0.5 μ l T4 DNA Ligase was added, and the reaction mix was incubated at 16°C overnight. 10 μ l dH₂O was added to the ligation reaction and 0.5 μ l was used to transform electro-competent *E. coli* DH10B cells. Typically, 6 to 10 colonies per μ l of transformed ligation mix were observed.

Example 12: Ligation of Activation Vectors to Genomic DNA and Transfection into Human Cells

Genomic DNA was harvested from a human cell line, HT1080 (10⁸ cells), according to published procedures (Sambrook et al., *Molecular Cloning*, Cold Spring Harbor Laboratory Press, (1989)). The isolated genomic DNA was digested with *Bam*HI under conditions that resulted in incomplete digestion. This

was accomplished by titrating the amount of *Bam*HI in the reaction. Each reaction contained 10 µg genomic DNA and *Bam*HI at a concentration of either 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.62, or 11.24 units. After a one hour incubation at 37°C, the reactions were stopped by phenol extraction, followed by ethanol precipitation. The digested DNA from each reaction was separated by agarose gel electrophoresis. Reactions containing DNA predominantly in the range of 10 kb to 400 kb were combined for ligation to the activation vector. The pooled, digested genomic DNA was then added to *Bam*HI linearized activation vector in 1X ligation buffer. Ligase (Life Technologies, Inc., 40 units) was added and the ligation reaction was incubated at 16°C for 24 hours. Following ligation, the genomic DNA/activation vector was transfected into HT1080 cells using LIPOFECTIN® (Life Technologies, Inc.) according to the manufacturer's procedures. Optionally, the HT1080 cells were irradiated prior to or after transfection. When cells were irradiated, doses in the range of 0.1 rads to 200 rads were found to be particularly useful. Following transfection, cells were grown in complete media. At 36 hours post-transfection, G418 (300 µg/ml) were added to the media. At 10-14 days post selection, the drug resistant clones were pooled, expanded, and harvested. Total RNA or mRNA was collected from the harvested cells. cDNA derived from vector activated genes was then synthesized and isolated using the methods described herein (see, e.g., Example 8 *supra*).

Example 13: Co-transfections of BAC Contig Clones with the Activation Vector

Genomic libraries were created in pUniBAC (Figure 34A-34B) according to published procedures (Shizuya et al., *Proc. Natl. Acad. Sci. USA* 89:8794 (1992)). Typically, the size of genomic fragments can be between 1 kb and 500 kb, and preferably between 50 kb and 500 kb. The BAC library was propagated in *E. coli*. To prepare plasmids for transfection, the library was plated onto LB agar plates containing 12.5 µg/ml chloramphenicol. Approximately 1000 clones were present on each 150 mm plate. Following growth and selection, the colonies

from each plate were eluted from the agar plate through the addition of LB and pooled. Each pool (~10,000 clones) was grown in 1 liter LB/12.5 µg/ml chloramphenicol overnight. BAC plasmids were then isolated from each pool using a commercial kit (Qiagen).

Purified BAC clones were digested with I-Ppo-I which cleaves a unique site in the BAC vector flanking the cloning site. Since I-Ppo-I is an ultra-rare cutter, it will not digest the vast majority of genomic DNA inserts. Following digestion, the linearized genomic library clones were cotransfected into HT1080 cells using LIPOFECTIN® (Life Technologies, Inc.) according to the manufacturer's directions. Briefly, 10 µg of BAC genomic DNA was combined with 1 µg of linearized pRIG20 (Figure 31A-31C) in α-MEM (no serum). 5 µg of LIPOFECTIN® was added to the DNA and the mixture was incubated at room temperature for 15 minutes. The DNA/LIPOFECTIN® mixture was then added to 10⁵ HT1080 cells in a 6 well dish. The cells were incubated with the DNA/LIPOFECTIN® in serum free α-MEM for 12 hours, washed, and placed in α-MEM/10%FBS for 36 hours. To select for cells that had integrated the vector and genomic DNA, the transfected cells were replated into a 10 cm dish and incubated in the presence of 300 µg/ml G418 for 10 days. Drug resistant clones were expanded and harvested to allow isolation of the activated cDNA molecules as described herein in Example 8.

Example 14: *In vitro* Integration of Activation Vector into Purified Genomic DNA and Transfection of the Integration Products into Host Cells

Genomic DNA was isolated and cloned into the Bacterial Artificial Chromosome, pUniBAC (Figure 34A-34B), using published procedures (Sambrook et al., *Molecular Cloning*, Cold Spring Harbor Laboratory Press, (1989); Shizuya et al., *Proc. Natl. Acad. Sci. USA* 89:8794 (1992)). Following ligation of the genomic inserts into pUniBAC, the plasmids were transformed into the *E. coli* strain DH10B (Life Technologies, Inc.) and selected on tetracycline.

Individual bacterial clones were combined into pools containing approximately 1000 members. Each pool was grown to saturation in 1 liter LB/tetracycline. pUniBAC plasmids containing genomic DNA inserts were isolated from the bacteria using a commercial kit (Qiagen).

For each pool of UniBAC clones, 2 µg of the library were incubated with 50 ng of the activation vector pRIG-T and 1 unit of mutant Tn5 transposase for 2 hours at 37°C (transposase available from Epicentre Technologies). Following incubation, the pUniBAC clones were transformed into DH10B cells and selected on chloramphenicol. All colonies from each pool were combined and grown in 1 liter LB/chloramphenicol. Plasmids were harvested using Qiagen Tip-500 columns according to the manufacturer's instructions.

For each pool, 20 µg of the library was transfected into 2×10^6 HT1080 cells with 30 µg Ex-gen 500 (MBI Fermentas) according to the manufacturer's instructions. At 48 hours post-transfection, the cells were placed into media containing 3 µg/ml puromycin. After 10 days of growth in the presence of puromycin, drug resistant clones were pooled, expanded and harvested for gene discovery. To isolate vector activated genes, mRNA from each pool of cells was isolated, converted to cDNA, and cloned into plasmids as described in Example 8. Individual cDNA clones were analyzed by restriction digestion and sequencing.

Example 15: Creation of Protein Expression Libraries from Cloned Genomic DNA

A genomic library containing genomic DNA inserts (100 kb avg. size) was created in pUniBAC as described in Examples 13 and 14. (Note: In some embodiments of the invention, the genomic fragments are cloned into the linearization site of an activation vector, wherein the activation vector is preferably a YAC, BAC, PAC, or Cosmid based vector.) In this example, the activation vector, pRIG-TP, was integrated into the BAC genomic library using in vitro transposition as described in Example 14. pRIG-TP is shown in Figure 36. Following integration, the library plasmids were transformed into E. coli and BAC

vectors containing an integrated pRIG-TP vector were selected for on chloramphenicol plates. Colonies were pooled and grown to saturation in LB/Tetracycline. BAC plasmids were harvested using a commercial kit (Qiagen).

For each transfection, 20 ug of the BAC library was transfected into 2×10^6 HT1080 cells using 30 ug Ex-gen 500 (MBI Fermentas) according to the manufacturer's instructions. At 48 hours post transfection, the cells were placed into media containing 3 ug/ml puromycin. After 10 days of selection, drug resistant clones were pooled and expanded. The expanded pools of drug resistant clones were divided into separate groups for freezing, protein production, and episome amplification.

To isolate and test activated secreted proteins, culture supernatants were harvested and saved at -80°C until used in specific assays. Activated intracellular proteins were harvested from cell lysates (prepared by any method known in the art) and used in in vitro assays.

To amplify the copy number of the BAC episomes, the cells were selected with increasing concentrations of methotrexate. In these experiments, the initial methotrexate concentration was 20 nM. Methotrexate concentrations were doubled every 7 days until cells resistant to 5 μM were obtained. At each methotrexate concentration, a portion of cells were removed for storage and protein production. Activated secreted and intracellular proteins were harvested from these cells as described for the non-methotrexate selected cells.

Having now fully described the present invention in some detail by way of illustration and example for purposes of clarity of understanding, it will be obvious to one of ordinary skill in the art that the same can be performed by modifying or changing the invention within a wide and equivalent range of conditions, formulations and other parameters without affecting the scope of the invention or any specific embodiment thereof, and that such modifications or changes are intended to be encompassed within the scope of the appended claims.

All publications, patents and patent applications mentioned in this specification are indicative of the level of skill of those skilled in the art to which this invention pertains, and are herein incorporated by reference to the same extent as if each individual publication, patent or patent application was specifically and individually indicated to be incorporated by reference.

5

WHAT IS CLAIMED IS:

1. A vector construct comprising:
 - (a) a first transcriptional regulatory sequence operably linked to a first unpaired splice donor sequence;
 - (b) a second transcriptional regulatory sequence operably linked to a second unpaired splice donor sequence; and
 - (c) a linearization site.

2. The vector construct of claim 1, wherein said linearization site is located between said first unpaired splice donor site and said second transcriptional regulatory sequence.

3. The vector construct of claim 1, wherein when said vector integrates into the genome of a host cell, said first transcriptional regulatory sequence is in an inverted orientation relative to the orientation of said second transcriptional regulatory sequence.

4. The vector of claim 1, wherein said vector has been rendered linear by cleavage at said linearization site.

5. A vector construct comprising, in sequential order:
 - (a) a transcriptional regulatory sequence;
 - (b) an unpaired splice donor site;
 - (c) a rare cutting restriction site; and
 - (d) a linearization site.

6. A vector construct comprising, in sequential order:
 - (a) a transcriptional regulatory sequence;
 - (b) a vector-encoded exon comprising a rare cutting restriction site;
 - (c) an unpaired splice-donor site; and

(d) a linearization site.

7. A vector construct comprising, in sequential order:

- (a) a transcriptional regulatory sequence;
- (b) a vector-encoded exon comprising a first rare cutting restriction site;
- (c) an unpaired splice-donor site;
- (d) a second rare cutting restriction site; and
- (e) a linearization site.

8. A vector construct comprising:

- (a) a first transcriptional regulatory sequence operably linked to a selectable marker lacking a polyadenylation signal; and
- (b) a second transcriptional regulatory sequence operably linked to an exon-splice donor site complex,

wherein said first transcriptional regulatory sequence is in the same orientation in said vector construct as said second transcriptional regulatory sequence.

9. A vector construct comprising a transcriptional regulatory sequence operably linked to a selectable marker lacking a polyadenylation signal, and further comprising an unpaired splice donor site.

10. A vector construct comprising a first transcriptional regulatory sequence operably linked to a selectable marker lacking a polyadenylation signal, and further comprising a second transcriptional regulatory sequence operably linked to an unpaired splice donor site.

11. The vector construct of any one of claims 1, 8, or 10, wherein said first transcriptional regulatory sequence or said second transcriptional regulatory sequence is a promoter.

12. The vector construct of claim 11, wherein said promoter is selected from the group consisting of a CMV immediate early gene promoter, an SV40 T antigen promoter, a tetracycline-inducible promoter, and a β -actin promoter.

13. The vector construct of any one of claims 5-7 or 9, wherein said transcriptional regulatory sequence is a promoter.

14. The vector construct of claim 13, wherein said promoter is selected from the group consisting of a CMV immediate early gene promoter, an SV40 T antigen promoter, a tetracycline-inducible promoter, and a β -actin promoter.

15. The vector construct of any one of claims 8-10, wherein said selectable marker is selected from the group consisting of a neomycin gene, a hypoxanthine phosphoribosyl transferase gene, a puromycin gene, a dihydroorotase gene, a glutamine synthetase gene, a histidine D gene, a carbamyl phosphate synthase gene, a dihydrofolate reductase gene, a multidrug resistance 1 gene, an aspartate transcarbamylase gene, a xanthine-guanine phosphoribosyl transferase gene, an adenosine deaminase gene, and a thymidine kinase gene.

16. A vector construct comprising:

- (a) a positive selectable marker,
- (b) a negative selectable marker; and
- (c) an unpaired splice donor site,

wherein said positive and negative selectable markers and said splice donor site are oriented in said vector construct in an orientation that results in expression of said positive selectable marker in active form, and either non-expression of said negative selectable marker or expression of said negative selectable marker in inactive form, when said vector construct is integrated into the genome of a eukaryotic host cell in such a way that an endogenous gene in said genome is activated.

17. The vector construct of claim 16, wherein said positive selection marker and said negative selection marker both lack a polyadenylation signal.

18. The vector construct of claim 16, wherein said positive selection marker is selected from the group consisting of a neomycin gene, a hypoxanthine phosphoribosyl transferase gene, a puromycin gene, a dihydroorotase gene, a glutamine synthetase gene, a histidine D gene, a carbamyl phosphate synthase gene, a dihydrofolate reductase gene, a multidrug resistance 1 gene, an aspartate transcarbamylase gene, a xanthine-guanine phosphoribosyl transferase gene, and an adenosine deaminase gene.

19. The vector construct of claim 16, wherein said negative selection marker is selected from the group consisting of a hypoxanthine phosphoribosyl transferase gene, a thymidine kinase gene, and a diphtheria toxin gene.

20. A eukaryotic host cell comprising the vector construct of any one of claims 1, 5-10, or 16.

21. The eukaryotic host cell of claim 20, wherein said cell is an animal cell.

22. The eukaryotic host cell of claim 21, wherein said animal cell is selected from the group consisting of a mammalian cell, an insect cell, an avian cell, an annelid cell, an amphibian cell, a reptilian cell, and a fish cell.

23. The eukaryotic host cell of claim 21, wherein said animal cell is a mammalian cell.

24. The eukaryotic host cell of claim 23, wherein said mammalian cell is a human cell.

25. The eukaryotic host cell of claim 20, wherein said cell is a plant cell.

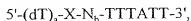
26. The eukaryotic host cell of claim 20, wherein said cell is a fungal cell.

5 27. The eukaryotic host cell of claim 26, wherein said fungal cell is a yeast cell.

28. The eukaryotic host cell of claim 21, wherein said cell is an isolated cell.

10 29. The eukaryotic host cell of claim 21, wherein said vector construct is integrated into the genome of said host cell.

30. A primer molecule comprising a PCR-amplifiable sequence and a degenerate 3' terminus, wherein said primer molecule has the structure:



15 wherein a is a whole number from 1 to 100, X is a PCR-amplifiable sequence consisting of a nucleic acid sequence of about 10-20 nucleotides in length, N is any nucleotide, and b is a whole number from 0 to 6.

31. The primer molecule of claim 30, wherein said PCR-amplifiable sequence comprises one or more restriction sites.

20 32. The primer molecule of claim 30, wherein a is a whole number from 10 to 30.

33. The primer molecule of claim 30, wherein said primer molecule comprises one or more hapten molecules conjugated to one or more bases of said primer molecule.

34. The primer molecule of claim 33, wherein said hapten molecules are selected from the group consisting of biotin, digoxigenin, an antibody, an enzyme, lipopolysaccharide, apotransferrin, ferrotansferrin, insulin, a cytokine an extracellular matrix protein, an integrin, ankyrin, C3bi, fibrinogen, spectrin, a
5 cytokine receptor, an insulin receptor, a transferrin receptor, polymyxin B, endotoxin-neutralizing protein (ENP), an enzyme-specific substrate, protein A, protein G, a cell-surface Fc receptor, an antibody-specific antigen, an antibody-specific peptide, avidin, and streptavidin.

35. The primer molecule of claim 33, wherein said hapten molecule is
10 biotin.

36. A method for first strand cDNA synthesis comprising:

- (a) annealing a primer of claim 30 to an RNA template molecule to form an primer-RNA complex; and
- (b) treating said primer-RNA complex with reverse transcriptase and one or more deoxynucleoside molecules under conditions favoring the reverse transcription of said primer-RNA complex to
15 synthesize a first strand cDNA.

37. A method for isolating an activated gene from a host cell genome, comprising:

- (a) introducing a vector comprising a transcriptional regulatory sequence, a vector-encoded exon, an unpaired splice donor site, and a vector-encoded intron into a host cell;
 - (b) allowing said vector to integrate into the genome of said host cell by non-homologous recombination, under conditions such that
20 said vector activates an endogenous gene in said genome;
 - (c) isolating RNA from said host cell;
 - (d) synthesizing first strand cDNA by reverse transcription of said isolated RNA;
- 25

- (e) annealing a primer specific for said vector-encoded exon to said first strand cDNA to create a primer-first strand cDNA complex; and
- (f) contacting said primer-first strand cDNA complex with a DNA polymerase under conditions favoring the production of a second strand cDNA product substantially complementary to said first strand cDNA.

38. A method for isolating an activated gene from a host cell genome, comprising:

- (a) introducing a vector comprising a transcriptional regulatory sequence, a vector-encoded exon, an unpaired splice donor site, and a vector-encoded intron into a plurality of host cells;
- (b) allowing said vector to integrate into the genomes of said host cells by non-homologous recombination, under conditions such that said vector activates an endogenous gene in said genomes;
- (c) cultivating said host cells under conditions favoring the production of a plurality of individual clones from said host cells, wherein each of said individual clones in said plurality of clones contains said vector integrated into a unique site in said host cell genome;
- (d) isolating RNA from said plurality of clones;
- (e) synthesizing first strand cDNA by reverse transcription of said isolated RNA;
- (f) annealing a first primer specific for said vector-encoded exon to said first strand cDNA to create a primer-first strand cDNA complex; and
- (g) contacting said primer-first strand cDNA complex with a DNA polymerase under conditions favoring the production of a second strand cDNA product substantially complementary to said first strand cDNA.

39. The method of claim 37, further comprising treating said second strand cDNA product with a restriction enzyme that cleaves at a restriction site located on said vector-encoded exon.

5 40. The method of claim 38, further comprising treating said second strand cDNA product with a restriction enzyme that cleaves at a restriction site located on said vector-encoded exon.

10 41. The method of claim 37, further comprising treating said second strand cDNA product with a restriction enzyme that cleaves at a restriction site located on said vector-encoded intron downstream of said unpaired splice donor site.

42. The method of claim 38, further comprising treating said second strand cDNA product with a restriction enzyme that cleaves at a restriction site located on said vector-encoded intron downstream of said unpaired splice donor site.

15 43. The method of claim 37, further comprising amplifying said second strand cDNA product using a second primer specific for said vector-encoded exon and a third primer specific for said first primer.

20 44. The method of claim 38, further comprising amplifying said second strand cDNA product using a second primer specific for said vector-encoded exon and a third primer specific for said first primer.

45. An isolated gene produced according to the method of any one of claims 37-44.

46. A host cell comprising the isolated gene of claim 45.

47. A vector comprising the isolated gene of claim 45.

48. The vector of claim 47, wherein said vector is an expression vector.

49. A method of producing a polypeptide, comprising:

- (a) introducing the vector of claim 47 into a host cell; and
(b) culturing said host cell under conditions favoring the expression by said host cell of a polypeptide encoded by said isolated gene.

50. The method of claim 49, further comprising isolating said polypeptide.

51. A polypeptide produced according to the method of claim 49 or claim 50.

52. A method of producing a polypeptide, comprising:

- (a) introducing into a host cell a vector comprising a transcriptional regulatory sequence operably linked to an exonic region followed by an unpaired splice donor site, under conditions favoring the integration of said vector into the genome of said host cell and resulting in the activation of an endogenous gene in said genome; and

- (b) culturing said host cell under conditions favoring the expression by said host cell of a polypeptide at least partially encoded by said exonic region,

wherein said exon contains a translational start site positioned at position -3, or at an increment of 3 bases upstream therefrom, from the 5'-most base of said splice donor site.

53. A method of producing a polypeptide, comprising:

- (a) introducing into a host cell a vector comprising a transcriptional regulatory sequence operably linked to an exonic region followed by an unpaired splice donor site, under conditions favoring the integration of said vector into the genome of said host cell and resulting in the activation of an endogenous gene in said genome; and
- (b) culturing said host cell under conditions favoring the expression by said host cell of a polypeptide at least partially encoded by said exonic region,

wherein said exon contains a translational start site positioned at position -2, or at an increment of 3 bases upstream therefrom, from the 5'-most base of said splice donor site.

54. A method of producing a polypeptide, comprising:

- (a) introducing into a host cell a vector comprising a transcriptional regulatory sequence operably linked to an exonic region followed by an unpaired splice donor site, under conditions favoring the integration of said vector into the genome of said host cell and resulting in the activation of an endogenous gene in said genome; and
- (b) culturing said host cell under conditions favoring the expression by said host cell of a polypeptide at least partially encoded by said exonic region,

wherein said exon contains a translational start site positioned at position -1, or at an increment of 3 bases upstream therefrom, from the 5'-most base of said splice donor site.

55. The method of any one of claims 52-54, further comprising isolating said polypeptide.

56. A polypeptide produced by any one of claims 52-54.
57. A polypeptide produced by the method of claim 55.

0081231-011800

Compositions and Methods for Non-targeted Activation of Endogenous Genes

Abstract

5 The present invention is directed generally to activating gene expression or causing over-expression of a gene by recombination methods *in situ*. The invention also is directed generally to methods for expressing an endogenous gene in a cell at levels higher than those normally found in the cell. In one embodiment of the invention, expression of an endogenous gene is activated or increased following integration into the cell, by non-homologous or illegitimate recombination, of a regulatory sequence that activates expression of the gene. In another embodiment, the expression of the endogenous gene may be further increased by co-integration of one or more amplifiable markers, and selecting for increased copies of the one or more amplifiable markers located on the integrated vector. In another embodiment, the invention is directed to activation of 10 endogenous genes by non-targeted integration of specialized activation vectors, which are provided by the invention, into the genome of a host cell. The invention also provides methods for the identification, activation, isolation, and/or expression of genes undiscoverable by current methods since no target sequence is necessary for integration. The invention also provides methods for isolation of 20 nucleic acid molecules (particularly cDNA molecules) encoding a variety of proteins, including transmembrane proteins, and for isolation of cells expressing such transmembrane proteins which may be heterologous transmembrane proteins. The invention also is directed to isolated genes, gene products, nucleic acid molecules, to compositions comprising such genes, gene products and nucleic acid molecules, to vectors and host cells comprising such genes and nucleic acid 25 molecules, that may be used in a variety of therapeutic and diagnostic applications. Thus, by the present invention, endogenous genes, including those associated with human disease and development, may be activated and isolated without prior knowledge of the sequence, structure, function, or expression profile of the genes. 30

09484331-011300

Random Activation of Gene Expression (RAGE)

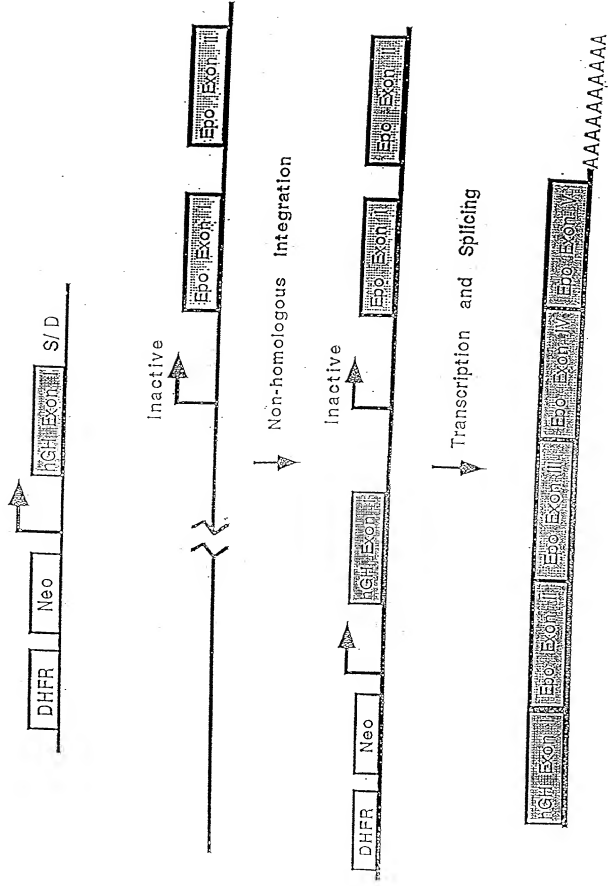
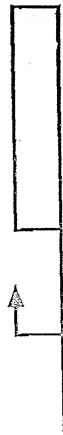


FIGURE 1

Activation Constructs without Translation Start Codons

Construct #



1



2



Untranslated

S/D

Splice Donor

Fig. 2

Construct #

002770-12EE+8+60

3-5



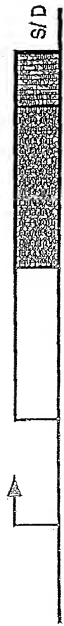
6-8



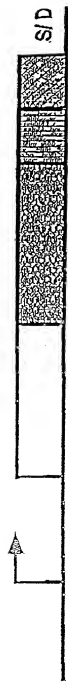
9-11



12-14



15-17



Untranslated



Secretion Signal



Protease Cleavage Site



Translated



Epitope Tag

S/D

Splice Donor

Fig. 3

pRIG-1

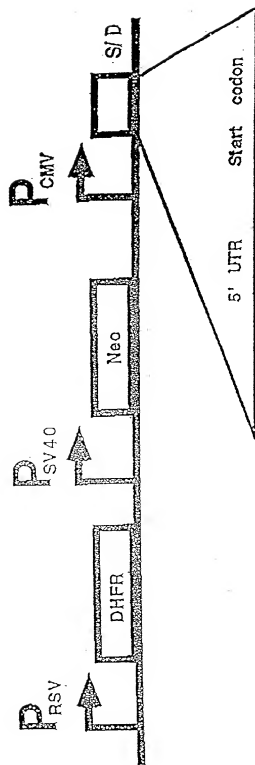


FIG. 4

5'AGATCTTCAATATTGGGCCATTAGCCATATTATTCATTGGTTATATAGCATAAATC
 AATATTGGCTATTGGCCATTGCATA
 CGTTGTATCTATATCATATAATATGTACATTTATATTGGCTCATGTCCAATATGACCG
 CCATGTTGGCATTGATTATTGACT
 AGTTATTAATAGTAATCAATTACGGGGTCATTAGTTCATAGCCCATATATGGAGT
 TCCGCGTTACATAAATTACGGTAA
 TGGCCCGCTGGCTGACCGCCCAACGACCCCCGCCATTGACGTCAATAATGAAG
 TATGTTCCCATAGTAACGCCAATAG
 GGACTTTCCATTGACGTCAATGGGTGGAGTATTACGGTAAACTGCCCACTTGGC
 AGTACATCAAGTGTATCATATGCCA
 AGTCCGCCCCCTATTGACGTCAATGAACGGTAAATGGCCCGCTGGCATTATGCC
 AGTACATGACCTTACGGGACTTTCC
 TACTTGGCAGTACATCTACGTATTAGTCATCGCTATTACCATGGTGATGCGGTTTT
 GGCAGTACACCAATGGGCGTGGAT
 AGCGGTTTGACTCAGGGGATTTCGAAGTCTCCACCCCATTTGACGTCAATGGGAG
 TTTGTTTGGGCACAAAATCAACGG
 GACTTTCCAAAATGTCTGAACAACCTGCGATGCGCCCGCCCGTTGACGCAATGGG
 CGGTAGGCGTGTACGGTGGGAGGTC
 TATATAAGCAGAGCTCGTTTTAGTGAACCGCTCAGATCACTAGAAGCTTTATTGCGG
 TAGTTTATCACAGTTAAATTGCTAA
 CGCAGTCAGTGCTTCTGACACAACAGTCTCGAACTTAAGTGCAGTGACTCTCTT
 AATTAACCTCCACCAGTCTCACTTCA
 GTTCCTTTTGCTCCACCAGTCTCACTTCAGTTCCCTTTGTCATGAAGAGCTCAGAA
 TCAAAAAGAGGAAACCAACCCCTAA
 GATGAGCTTTCCATGTAAATTTGTAGCCAGCTTCTCTGATTTTCAATGTTTCTT
 CCAAAGGTGCAGTCTCCAAAGAGA
 TTACGAATGCCTTGGAACCTGGGGTGCCTTGGGTGAGGACATCAACTTGGACAT
 TCCTAGTTTTCAAATGAGTGATGAT
 ATTGACGATATAAATGGGAAAAAATTCAGACAAGAAAAAGATTGCACAATTCA
 GAAAAGAGAAAAGAGACTTTCAAGGA
 AAAAGATACATATAAGCTATTAAAAATGGAACTCTGAAAATTAAGCATCTGAAG
 ACCGATGATCAGGATATCTACAAGG
 TATCAATATATGATACAAAAGGAAAAAATGTGTTGGAAAAAATATTGATTTGAA
 GATTCAAGAGAGGGTCTCAAAACCA
 AAGATCTCTGGACTTGTATCAACACAACCTGACCTGTGAGGTAATGAATGGAA
 CTGACCCCGAATTAAACCTGTATCA
 AGATGGGAAACATCTAAAACTTTCTCAGAGGGTTCATCACACAAGTGGACCACC
 AGCCTGAGTGCAAAAATCAAGTGCA
 CAGCAGGAAACAAAGTCAGCAAGGAATCCAGTGTGCGACCTGTGACGTGTCCAG
 AGAAAGGGATCCAGGTGAGTAGGGCC
 CGATCCTTCTAGAGTCGAGCTCTCTTAAGGTAGCAAGGTTACAAGACAGGTTTAA
 GGAGACCAATAGAACTGGGCTTGT
 CGAGACAGAGAAGACTTTGCGTTTCTGATAGGCACCTATTGGTCTTACGCGGCC
 GCGAATTCCAAGCTTGAGTATTCTA
 TCGTGTACCTAAATAACTTGGCGTAATCATGGTCTATCTGTTTCTGTGTGAA
 ATTGTTATCCGCTCACAATTCACA
 CAACATACGAGCCGGAAGCATAAAGTGTAAAGCCTGGGGTGCCTAATGAGTGA
 CTAACCTCACATTAATTGCGTTGCGCGATGCTTCCATTTTGTGAGGGTTAATGC-

Figure 5A

TTCGAGAAGACATGATAAGATACATTGATGAGTTTGGACAAACCAACAAGAAT
 GCAGTGA AAAAATGCTTTATTTGTGAAATTTGTGATGCTATTGCTTTATTGTAA
 CCATTATAAGCTGCAATAAACA
 AGTTAACAAACAACAATTGCATTCTTTATGTTTCAGGTTTCAGGGGGAGATGTGG
 GAGGTTTTTTAAAGCAAGTAAACCC
 TCTACAAATGTGTTAAATCCGATAAGGATCGATTCCGGAGCCTGAATGGCGAAT
 GGACGCGCCCTGTAGCGGCGCATT
 AGCGCGGCGGGTGTGGTGGTTACGCGCACGTGACCGCTACACTTGCCAGCGCCC
 TAGCGCCCGCTCTTTTCGCTTTCTTC
 CCTTCCTTTCTCGCCACGTTTCGCGGCTTCCCGTCAAGCTCTAAATCGGGGGC
 TCCTTTTAGGGTTCCGATTAGTGC
 TTTACGGCACCTCGACCCCAAAAAACFTGATTAGGGTGATGGTTACGTTAGTGGG
 CCATCGCCCTGATAGACGGTTTTTC
 GCCCTTTGACGTTGGGAGTCCACGTTCTTTAATAGTGGACTCTTGTTCCAAACTGG
 AACAACACTCAACCCCTATCTCGGTC
 TATTCTTTTGATTTATAAGGGATTTTGGCGATTTTCGGCTATTGGTTAAAAAATGA
 GCTGATTTAACAAAAATTAAACGC
 GAATTTTAACAAAAATTAACGCTTACAATTTCCGCTGTGTACCTTCTGAGGCGG
 AAAGAACCAGCTGTGGAATGTGTGT
 CAGTTAGGGTGTGGAAGTCCCCAGGCTCCCCAGCAGGCAGAAAGTATGCAAAGC
 ATGCATCTCAATTAGTCAGCAACCAG
 GTGTGGAAGTCCCAGGCTCCCCAGCAGGCAGAAAGTATGCAAAGCATGCATCT
 CAATTAGTCAGCAACCATAGTCCCGC
 CCTTAACCTCCGCCCATCCCGCCCTAACTCCGCCAGTTCCGCCCATCTCCGCC
 CCAATGGCTGACTAATTTTTTTTATT
 TATGACAGAGGCCGAGGCCCTCGGCTCTGAGCTATTCCAGAAGTAGTGAGGA
 GGCTTTTTTTGGAGGCCTAGGCTTTTG
 CAAAAAGCTTGATTCTCTGACACAACAGTCTCGAACTTAAGGCTAGAGCCACCA
 TGATTGAACAAGATGGATTGCACGC
 AGGTTCTCCGGCCGCTGGGGTGGAGAGGCTATTCCGCTATGACTGGGCACAACAG
 ACAATCGGCTGTCTGTATGCCGCCG
 TGTTCGGCTGTGAGCGCAGGGGCCCGGTTCTTTTTGTCAAGACCGACCTGTC
 CGGTGCCCTGAATGAACCTGCAGGAC
 GAGGCAGCGCGGCTATCGTGGCTGGCCACGACGGGCGTTCTTGCGCAGCTGTG
 CTCGACGTTGTCACTGAAGCGGGAAG
 GGACTGGCTGCTATTGGGCGAAGTGCCGGGGCAGGATCTCTGTCTATCTCACCTT
 GCTCCTGCGGAGAAAGTATCCATCA
 TGGCTGATGCAATGCGCGGCTGCATACGCTGTATCCGGCTACCTGCCCATTCGA
 CCACCAAGCGAAACATCGCATCGAG
 CGAGCACGTACTCGGATGGAAGCCGGTCTTGTGATCAGGATGATCTGGACGAA
 GAGCATCAGGGGCTCGGCCAGCCGA
 ACTGTTTCGCCAGGCTCAAGGCGCGCATGCCGACGGCGAGGATCTCGTCTGTGAC
 CCATGGCGATGCGCTGTGCGGAATA
 TCATGTGTGAAAAATGGCCGCTTTTCTGGATTTCATCGACTGTGGCCGGCTGGGTGT
 GGGCGACCGCTATCAGGACATAGCG
 TTGGCTACCCGCTGATATTGCTGAAGAGCTTGGCGCGAATGGGCTGACCGCTTCC
 TCGTGTCTTACGCTATCGCGCTCC
 CGATTGCGACGCGCATCGCCTTCTATCGCCTTCTTGACGAGTCTTCTGAGCGGGA
 CTCTGGGGTTCGAAATGACCGACCAAGCGACGCCCAACCTGCCATCAGATGGC-

Figure 5B

CGCAATAAAATATCITTTATTTTCATTACATCTGTGTGTGGTTTTTTTGTGTGAAGA.
TCCGCGTA-
TGGTGCACCTCTCAGTACAAATCTGCTCTGATGCCGCATAGTTAAGCCAGCCCCGAC
ACCCGCCAACAC
CCGCTGACGCGCCCTGACGGGCTTGTCTGCTCCCGGCATCCGCTTACAGACAAGC
TGTGACCGTCTCCGGGAGCTGCATG
TGTGAGAGGTTTTTACCGTTCATCACCAGAAACGCGGAGACGAAAGGGCCTCGTGA
TACGCCCTATTTTTATAGGTTAATGT
CATGATAATAATGGTTTTCTTAGACGTCAGGTGGCACTTTTCGGGGAAATGTGCGC
GGAACCCCTATTTGTTTTATTTTCT
AAATACATTCAAATATGTATCCGCTCATGAGACAATAACCTGATAAATGCTTCA
ATAATATTGAAAAGGAAGAGTATG
AGTATTCAACATTTCCGTGTGCGCCCTTATCCCTTTTTTTCGGGCATTTTGCTTCC.
TGTTTTTGCTCAGCCAGAAACGCT
GGTGAAGTAAAGATGCTGAAGATCAGTTGGGTGCACGAGTGGGTTACATCGA
ACTGGATCTCAACAGCGGTAAGATCC
TTGAGAGTTTTTCGCCCGAAGAACGTTTTTCCAATGATGAGCACITTTAAAGTTCT
GCTATGTGGCGCGGTATTATCCCGT
ATTGACGCGGGCAAGCAACTCGGTGCGCCCATACACTATTCTCAGAATGACT
TGGTTGAGTACTCACAGTCAAGCA
AAAGCATCTTACGGATGGCATGACAGTAAGAGAATTATGCAGTGTGCCATAACC
ATGAGTGATAACACTGCGGCCAACT
TACTTCTGACAAACGATCGGAGGACCGAAGGAGCTAACCGCTTTTTTGACAAACAT
GGGGGATCATGTAACTCGCCTTGAT
CGTTGGGAACCGGAGCTGAATGAAGCCATACCAAACGACGAGCGTGACACCACG
ATGCCTGTAGCAATGGCAACAACGTT
GCGCAAACTATTAACTGGCGAACTACTTACTCTAGCTTCCCGGCAACAATTAATA
GACTGGATGGAGGCGGATAAAGTTG
CAGGACCACTTCTGCGCTCGGCCCTTCCGGCTGGCTGGTTTATTGCTGATAAAATC
TGGAGCCGGTGAGCGTGGGTCTCGC
GGTATCATTTGAGCACTGGGGCCAGATGGTAAGCCCTCCCGTATCGTAGTTATCT
ACACGACGGGGAGTCAGGCAACTAT
GGATGAACGAAATAGACAGATCGCTGAGATAGGTGCCTCACTGATTAAGCATTGG
TAACTGTGACACCAAGTTTACTCAT
ATATACTTTAGATTGATTTTAAAACTTCACTTTTAAATTTAAAGGATCTAGGTGAAG
ATCCTTTTTGATAATCTCATGACC
AAAATCCCTTAACTGAGTTTTTCGTTCCACTGAGCGTCAGACCCCGTAGAAAAGA
TCAAAGGATCTTCTTGAGATCCCTT
TTTTCTGCGCGTAATCTGCTGCTTGCAACAAAAAACACCGCTACCAGCGGTG
GTTTGTGTCGGGATCAAGAGCTAC
CAACTCTTTTTCCGAAGGTAACTGGCTTCAGCAGAGCGCAGATACCAAATACTGT
CCTTCTAGTGTAGCCGTAGTTAGGC
CACCCTTCAAGAACTCTGTAGCACCGCTACATACCTCGCTCTGCTAATCCTGT
TACCAGTGGCTGCTGCCAGTGGCGA
TAAGTCGTGTCTTACCGGGTTGGACTCAAGACGATAGTTACCGGATAAGGCGCAG
CGGTGCGGCTGAACGGGGGGTTGCT
GCACACAGCCAGCTTGGAGCGAACGACCTACACCGAACTGAGATACCTACAGC
GTGAGCTATGAGAAAGCGCCACGCTT
CCCGAAGGGAGAAAGCGGACAGGTATCCGGTAAGCGGCAGGGTCCGAACAGG-

Figure 5C

AGAGCGCACGAGGGAGGTTCCAGGGGGAAACGCCTGGTATCTTTATAGTCCTGTC
GGGTTTCGCCACCTCTGACTTGAGCGTCGATTTTGTGATGCTCGTCAGGGG
GGCGGAGCCTATGAAAAACGCCAGCAACGCGGCCCTTTTACGGTTTCCTGGCCTT
TTGCTGGCCTTTTGCTCACATGGCT
CGAC3'

0044331.01800

Figure 5D

5'AGATCTTCAATATTGGCCATTAGCCATATTATTCATTGGTTATATAGCATAAATC
 AATATTGGCTATTGGCCATTGGCAT
 ACGTTGTATCTATATCATAATATGTACATTTATATTGGCTCATGTCCAATATGACC
 GCCATGTTGGCATTGATTATTGAC
 TAGTTATTAAATAGTAATCAATTACGGGGTCATTAGTTCATAGCCCATATATGGAG
 TTCCGCGTTACATAACITACGGTAA
 ATGGCCCGCCTGGCTGACCGCCCAACGACCCCCGCCCATTGACGTCAATAATGAC
 GTATGTTCCCATAGTAAACGCCAATA
 GGGACTTTCCATTGACGTCAATGGGTGGAGTATTTACGGTAAACTGCCCACTGG
 CAGTACATCAAGTGTATCATATGCC
 AAGTCCGCCCCCTATTGACGTCAATGACGGTAAATGGCCCGCCTGGCATTATGCC
 CAGTACATGACCTTACGGGACTTTC
 CTACTTGGCAGTACATCTACGTATTAGTCATCGCTATTACCATGGTGTATGGGTT
 TTGGCAGTACACCAATGGGCGTGGA
 TAGCGGTTTGACTCACGGGGATTTCGAAGTCTCCACCCCATTGACGTCAATGGGA
 GTTTGTTTGGCACCAAAATCAACG
 GGACTTTCCAAAATGTCTGAACAACTCGCATCGCCCGCCCGTTGACGCAAAATGG
 GCGGTAGGCGTGTACGGTGGGAGGT
 CTATATAAGCAGAGGCTCGTTTAGTGAACCGTCAGATCACTAGAAGCTTTATTGCG
 GTAGTTTATCACAGTTAAATTGCTA
 ACGCAGTCAGTGCTTCTGACACAACAGTCTCGAACTTAAGCTGCAGTGACTCTCT
 TAATTAACCTCCACAGTCTCACTTC
 AGTTCCTTTTGCTCCACAGTCTCACTTCAGTTCCTTTTGCATGAAGAGCTCAGA
 ATCAAAAGAGGAAACCAACCCCTA
 AGATGAGCTTTCCATGTAATTTGTAGCCAGCTTCCTTCTGATTTTCAATGTTTCT
 TCCAAAGGTGCAGTCTCCAAAGAG
 ATTACGAATGCCTTGGAAACCTGGGGTGCCTTGGGTGAGGACATCAACTTGGACA
 TTCCTAGTTTTCAAATGAGTGATGA
 TATTGACGATATAAAATGGGAAAAAACTTCAGACAAAGAAAAAGATTGCACAATTC
 AGAAAAAGAGAAAGAGACTTTCAAGG
 AAAAAAGATACATATAAGCTATTTAAAAATGGAACTCTGAAAAATTAAGCATCTGAA
 GACCGATGATCAGGATATCTACAAG
 GTATCAATATATGATACAAAAGGAAAAAATGTGTTGGAAAAAATATTGATTTGA
 AGATTCAAGAGAGGGTCTCAAAACC
 AAAGATCTCTGGACTTGTATCAACACAACCTGACCTGTGAGGTAATGAATGGA
 ACTGACCCCGAATTAAACCTGTATC
 AAGATGGGAAACATCTAAACTTTCTCAGAGGGTCATCACACAAAGTGGACCAC
 CAGCCTGAGTGCAAAATTCAAGTGC
 ACAGCAGGGAAACAAAGTCAGCAAGGAATCCAGTGTGAGCCTGTCACTGTGCA
 GAGAAAGGGATCCCAGGTGAGTAGGG
 CCCGATCCTTCTAGAGTCGAGCTCTCTTAAGGTAGCAAGGTTACAAGACAGGTTT
 AAGGAGACCAATAGAAACTGGGCTT
 GTCGAGACAGAGAAGACTCTTGCGTTTCTGATAGGCACCTATTGGTCTTACGCGG
 CCGCGAATCCAAGCTTGAGTATTC
 TATCGTGTACCTAAATAACTTGGCGTAATCATGGTTCATATCTGTTTCTGTGTGA
 AATTGTTATCCGCTCACAAATTCGA
 CACAACATACGAGCCGGAAGCATAAAGGTAAAGCCTGGGGTGCCATATGAGTG
 AGCTAACTCACATTAATTGCGTTGCG
 CGATGCTTCATTTTGTGAGGGTTAATGCTTCGAGAAGACATGATAAGATACATT
 GATGAGTTTGGACAAACCACAACAAGAAATGCAGTGAAAAAATGCTTTATTGT-

Figure 6A

GAAATTTGTGATGCTATTGCTTTATTTGTAAACCATTATAAGCTGCAATAAA
 CAAGTTAAACAACAACAATTGCATTTCATTTATGTTCAGGTTTCAGGGGGAGATGT
 GGGAGGTTTTTTAAAGCAAGTAAAA
 CCTCTACAAATGTGGTAAATCCGATAAGGATCGATTCCGGAGCCTGAATGGCGA
 ATGGACGGCGCCTGTAGCGGCGCAT
 TAAGCGCGCGGGGTGTGGTGTTACCGCGCACGTGACCGCTACACTTGCCACGGC
 CCTAGCGCCCGCTCCTTTTCGCTTCT
 TCCCTTCCTTTCTCGCCACGTTGCGCGGCTTTCCCGTCAAGCTCTAAATCGGGG
 GCTCCCTTTAGGGTTCCGATTIAGT
 GCTTTACGGCACCTCGACCCCAAAAACCTGATTAGGGTGATGGTTACGTTAGTG
 GGCCATCGCCCTGATAGACGGTTTT
 TCGCCCTTTGACGTTGGAGTCCACGTTCTTTAATAGTGGACTCTTGTTCCAAACCTG
 GAACAACACTCAACCTATCTCGG
 TCTATTCTTTGATTTATAAGGGATTITGCGGATTTCGGCCTATTGGTTAAAAAAT
 GAGCTGATTAAACAAAAATTAAC
 GCGAATTTTAAACAAAAATTAACGCTTACAATTTGCGCTGTGTACCTTCTGAGGGC
 GGAAGAAGAACAGCTGTGGAATGTGT
 GTCAGTTAGGGGTGTGGAAAGTCCCGAGGCTCCCGAGCAGGCAGAGTATGCAAA
 GCATGCATCTCAATTAGTCAGCAACC
 AGGTGTGGAAGTCCCGAGGCTCCCGAGCAGGCAGAGTATGCAAAGCATGCAT
 CTCAATTAGTCAGCAACCATAGTCCC
 GCCCTAACTCCGCCATCCCGCCCTAACTCGGCCAGTTCGGCCCATCTCCG
 CCCATGGCTGACTAATTTTTTTA
 TTTATGCAGAGGCGAGGCGCGCTCGGCCCTGTGAGCTATTCCAGAAGTAGTGAGG
 AGGCTTTTTTGGAGGCGTAGGCTTT
 TGCAAAAAGCTTGATTCTTCTGACACAACAGTCTCGAACTTAAGGCTAGAGCCAC
 CATGATTGAACAAGATGGATTGCAC
 GCAGGTTCTCCGGCGGCTTGGGTGGAGAGGCTATTCCGGCTATGACTGGGCACAA
 AGACAATCGGCTGCTCTGATGCCGC
 CGTGTTCGGGCTGTGAGCGCAGGGGCGCCCGGTTCTTTTTGTCAAGACCGACCTG
 TCCGGTGCCTGAATGAACTGCAGG
 ACGAGGCAGCGCGGCTATCGTGGCTGGCCACGACGGGCGTTCTTGTGCGCAGCTG
 TGCTCGACGTTGTCACTGAAGCGGGA
 AGGGACTGGCTGCTATTGGGCGAAGTGCCGGGGCAGGATCTCCTGTCTATCTACC
 TTGCTCCTGCGGAGAAAGTATCCAT
 CATGGCTGATGCAATGCGGCGGCTGCATACGCTTGATCCGGCTACCTGCCCATTC
 GACCAACCAAGCGAAACATCGCATCG
 AGCGAGCAGCTACTCGGATGGAAGCCGGTCTTGTGATCAGGATGATCTGGACG
 AAGAGCATCAGGGGCTCGCGCCAGCC
 GAACGTGTCGCCAGGCTCAAGGCGCGCATGCCGACGGCGAGGATCTCGTCTGTG
 ACCCATGGCGATGCCTGCTTGCCGAA
 TATCATGGTGGAAAAATGGCGCTTTTCTGGATTATCGACTGTGGCCGGCTGGGT
 GTGGCGGACCGCTATCAGGACATAGCGTTGGCTACCCGTGATATTGCTGAAGAGC
 TTGGCGGCGAATGGGCTGACCGCTTCTCGTGCTTTACGATATCGCCGCT
 CCCGATTGCGAGCGCATCGCCTTCTATCGCCTTCTTGACGAGTCTTCTGAGCGG
 GACTCTGGGGTTTCAAAATGACCGAC
 CAAGCGACGCCAACCTGCCATCAGATGGCCGCAATAAAATATCTTTATTTTCA
 TTACATCTGTGTGTGGTTTTTGT
 GTGAAGATCCGCGTATGGTGCACTCTCAGTACAATCTGCTCTGATGCCGATAGT
 TAAGCCAGCCCCGACACCCGCCAACCCCGCTGACGCGCCTGACGGGCT

Figure 6B

TGTCGCTCCCGGCATCCGCTTACAGACAAGCTGTGACCGTCTCCGGGAGCTGCA
 TGTGTCAGAGGTTTTCACCGTCATCACGAAACGCGGAGACGAAAGGGCTCGT
 GATACGCCTATTTTATAGGTTAAAT
 GTCATGATAATAATGGTTTCTTAGACGTCAGGTGGCACTTTTCGGGGAAATGTGC
 GCGGAACCCCTATTGTGTTATTTTT
 CTAATAACATTCAAATATGTATCCGCTCATGAGACAATAACCCCTGATAAATGCCT
 CAATAATATTGAAAAAGGAAGAGTA
 TGAGTATTCAACATTTCGCTGTGCGCCCTATTCCCTTTTTTGCGGCATTTTGCCCT
 CCTGTTTTTGCTCACCCAGAAACG
 CTGGTGAAAGTAAAGATGCTGAAGATCAGTTGGGTGCACGAGTGGGTTACATC
 GAACTGGATCTCAACAGCGGTAAGAT
 CCTTGAGAGTTTTTCGCCCCGAAGAACGTTTTCCAATGATGAGCACTTTTAAAGTT
 CTGCTATGTGGCGCGGTATTATCCC
 GTATTGACCGCGGCAAGAGCAACTCGGTGCGCCGCATACACTATTCTCAGAATGA
 CTTGGTTTGAGTACTCACCAGTCACA
 GAAAAGCATCTTACGGATGGCATGACAGTAAGAGAATTATGCAGTGTGCGCTAA
 CCATGAGTGATAACACTGACGCGCCAA
 CTTACTTCTGACAAAGATCGGAGGACCGAAGGAGCTAACCGCTTTTTTGACACAAC
 ATGGGGGATCATGTAACCTCGCCTTG
 ATCGTTGGGAACCGGAGCTGAATGAAGCCATACCAAACGACGAGCGTGACACCA
 CGATGCCTGTAGCAATGGCAACAACG
 TTGCGCAAACCTATTAAGTGGCGAACTACTTACTCTAGCTTCCCGGCAACAATTAA
 TAGACTGGATGAGGCGGATAAAGT
 TGCAAGACCACCTTCTGCGCTCGGCCCTTCCGGCTGGCTGGTTTATTGCTGATAAA
 TCTGGAGCCGGTGAGCGTGGGTCTC
 GCGGTATCATTTGACGACTGGGGCCAGATGGTAAGCCCTCCCGTATCGTAGTTAT
 CTACACGACGGGAGTCAAGGCAACT
 ATGGATGAACGAAATAGACAGATCGCTGAGATAGGTGCCTCACTGATTAAAGCATT
 GGTAACGTGCAGACCAAGTTTACTC
 ATATATACTTTAGATTGATTAAAAACCTCATTTTTTAATTTAAAGGATCTAGGTGA
 AGATCCTTTTTGATAATCTCATGA
 CCAAAATCCCTTAACGTGAGTTTTCGTTCCACTGAGCGTCAGACCCCGTAGAAAA
 GATCAAAGGATCTTCTTGAGATCCT
 TTTTTTCTGCGCTAATCTGCTGCTTGCAAAACAAAAAACCCGCTACCGCGG
 TGGTTTGTGTCGGGATCAAGAGCT
 ACCAACTCTTTTTCCGAAGTAACTGGCTTCAGCAGAGCGCAGATACCAAATACT
 GTCCTTCTAGTGATGCCGTAGTTAG
 GCCACCACCTCAAGAACTCTGTAGCACCGCCTACATACCTCGCTCTGCTAATCCT
 GTTACCAAGTGGCTGCTGCCAGTGGCGATAAGTCGTGCTTACCGGGTTGGACTCA
 AGACGATAGTTACCGGATAAGGCGCAGCGTGGGCTGAACGGGGGGTTT
 GTGCACACAGCCAGCTTGGAGCGAACGACCTACACCGAACTGAGATACCTACA
 GCGTGAGCTATGAGAAAGCGCCACGC
 TTCCCGAAGGGAGAAAGCGCGACAGGTATCCGGTAAGCGGCAGGGTCCGGAACAG
 GAGAGCGCACGAGGGAGCTTCCAGGG
 GGAACGCGCTGTGATCTTTATAGTCTGTGGGTTTCGCCACCTCTGACTTGAGC
 GTCGATTTTTGTGATGCTCGTCAGG
 GGGCGGAGCCCTATGAAAAACGCCAGCAACGCGGCTTTTTACGGTTCTTGGC
 CTTTTGCTGGCCTTTTGCTCACATGG
 CTCGAC3'

Figure 6C

5'AGATCTTCAATATTGGCCATTAGCCATATTATTCATTGGTTATATAGCATAAATC
 AATATTGGCTATTGGCCATTGCAT
 ACGTTGTATCTATATCATAATATGTACATTATATTGGCTCATGTCCAATATGACC
 GCCATTGTTGCCATTGATTATTGAC
 TAGTTATTAATAGTAATCAATTACGGGGTCATTAGTTCATAGCCCATATATGGAG
 TTCCGCGTTACATAAATTACGGTAA
 ATGGCCCGCCTGGCTGACCGCCCAACGACCCCGCCCATGACGTCAATAATGAC
 GTATGTTCCCATAGTAACGCCAATA
 GGGACTTTCCATTGACGTCAATGGGTGGAGTATTTACGTTAACTGCCCACTTGG
 CAGTACATCAAGTTGATCATATGCC
 AAGTCCGCCCCCTATTGACGTCAATGACGCTAAATGGCCCCGCTGGCATTATGCC
 CAGTACATGACCTTACGGGACITTC
 CTACTTGGCAGTACATCTACGTATTAGTCATCGCTATTACCATGGTGATGCGGTT
 TTGGCAGTACACCAATGGGCGTTGGA
 TAGCGGTTTGACTCACGGGGATTTTCCAAGTCTCCACCCCATGACGTCAATGGGA
 GTTTGTTTGGGCACCAAAATCAACG
 GGACTTTCCAAAATGTGTAACAACTGCGATCGCCCGCCCGTTGACGCCAAATGG
 GCGGTAGGCGGTGACGGTGGGAGGT
 CTATATAAGCAGAGCTCGTTTAGTGAAACCGTCAGATCACTAGAAGCTTTATTGCG
 GTAGTTTATCACAGTTAAATTGCTA
 ACGCAGTCAGTGTCTCTGACACAACAGTCTCGAACTTAAGCTGCAGTGACTCTCT
 TAATTAACTCCACCAGTCTCACTTC
 AGTTCCTTTTGGCTCCACCACTCACTTCAGTTCCTTTTGCATGAAGAGCTCAGA
 ATCAAAAGAGGAAACCAACCCCTA
 AGATGAGCTTTCCATGTAAATTTGTAGCCAGCTTCCTTCTGATTTTCAATGTTTCT
 TCCAAAGGTGCAGTCTCCAAAGAG
 ATTACGAATGCCTTGGAAACCTGGGGTGCCTTGGGTGACGACATCAACTTGGACA
 TTCCTAGTTTTCAAATGAGTGATGA
 TATTGACGATATAAAATGGGAAAAAACCTCAGACAAGAAAAAGATTGCACAATTC
 AGAAAAGAGAAAAGAGACTTTCAAGG
 AAAAAAGATACATATAAGCTATTTAAAAATGGAACCTTGAAAAATTAAGCATCTGAA
 GACCGATGATCAGGATATCTACAAG
 GTATCAATATATGATACAAAAGGAAAAAATGTGTTGGAAAAAATATTGATTGTA
 AGATTCAAGAGAGGGTCTCAAAACC
 AAAGATCTCCTGGACTTGTATCAACACAACCCCTGACCTGTGAGGTAATGAATGGA
 ACTGACCCCGAATTAAACCTGTATC
 AAGATGGGAAACATCTAAAAACTTTCTCAGAGGGTCATCACACAAAGTGGACCAC
 CAGCCTGAGTGCAAAATTCAGTGC
 ACAGCAGGGAACAAAGTCAGCAAGGAATCCAGTGTGAGCCTGTGAGCTGTCCA
 GAGAAAGGGATCCACAGGTGAGTAGG
 GCCCGATCCTTCTAGAGTCGAGCTCTTAAGGTAGCAAGGTTACAAGACAGGTT
 TAAGGAGACCAATAGAAACTGGGCT
 TGTGAGACAGAGAAGACTCTTGCGTTTCTGATAGGCACTTATGGTCTTACGCG
 GCCGCGAATTCGAAGCTTGAATTT
 CTATCGTGTACCTAAATAAATACTTGGCGTAATCATGGTCATATCTGTTTCTGTGTG
 AAATTGTTATCCGCTCACAATTC
 ACACAACATACGAGCCGGAAGCATAAAGTGTAAGCCTGGGGTGCCTAATGAGT
 GAGCTAACTCACATTAATTGCGTTGC
 GCGATGCTTCCATTTTGTGAGGGTTAATGCTTCGAGAAGACATGATAAGATACAT
 TGATGAGTTTGGACAAACCACAACA AGAATGCAGTGAAAAAATGC-

Figure 7A

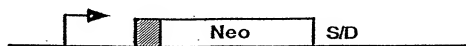
TTTATTTGTGAAATTGTGATG
 CTATTGCTTTATTTGTAACCAATTATAAGCTGCAATAA
 ACAAGTTAAACAACAATTGCAATTCATTTTATGTTTTCAGGTTTCAGGGGGAGATG
 TGGGAGGTTTTTTAAAGCAAAGTAA
 ACCTCTACAAATGTGGTAAATCCGATAAGGATCGATTCCGAGAGCTGAATGGCG
 AATGGACGCGCCCTGTAGCGGCGCA
 TTAAGCGCGCGGGTGTGGTGTGTTAAGCGCAAGTGACCGCTACACTTGCCAGCGC
 CCTAGCGCCGCTCCTTTGCGCTTC
 TTCCCTTCCTTTCTCGCCACGTTTCGCGGCTTTCCCGCTCAAGCTCTAAATCGGGG
 GCTCCCTTTAGGGTTCCGATTTAG
 TGCCTTACGGCACCTCGACCCCAAACTTGATTAGGGTGTAGGTTTACAGTAGT
 GGGCCATCGCCCTGATAGACGGTTT
 TTGCGCCCTTTGACGTTGGAGTCCACGTTCTTTAATAGTGGACTCTTGTTCCAAACT
 GGAACAACACTCAACCTATGTCG
 GTCTATTCTTTTGATTTATAAGGGATTTTGCCGATTTTCGGCTATTGGTTAAAAA
 TGAGCTGATTTTAAACAAAAATTTAA
 CGCGAATTTTAAACAAAAATTTAAACGCTTACAATTTTCGCTGTGACCTTCTGAGG
 CGGAAAGAACCAAGCTGTGAATGTG
 GTGCAGTTAGGGTGTGGAAAGTCCCCAGGCTCCCCAGCAGGCAGAAGTATGCAA
 AGCATGCATCTCAATTAGTCAGCAAC
 CAGGTGTGGAAAGTCCCAAGGCTCCCCAGCAGGCAGAAGTATGCAAAGCATGCA
 TCTCAATTAGTCAGCAACCATAGTCC
 CGCCCTAACTCCGCCATCCCGCCCTAACTCCGCCAGTTCCGCCCATTTCTCC
 GCGCCATGGCTGACTAATTTTTTTT
 ATTTATGCAGAGGCGGAGGCGCCCTCGGCTCTGAGCTATTCCAGAAAGTAGTGAG
 GAGGCTTTTTGAGGCCCTAGGCTT
 TTGCAAAAAAGCTTGATTTCTTGACACAACAGTCTCGAACTTAAGGCTAGAGCCA
 CCATGATTGAACAAGATGGATTGCA
 CGCAGGTTCTCCGGCCGCTTGGGTGGAGAGGCTATTCCGGCTATGACTGGGCACAA
 CAGACAATCGGCTGCTCTGATGCCG
 CCGTGTTCGGCGTGTGACGCGCAGGGGCGCCCGGTTCTTTTGTCAAGACCGACCT
 GTCCGGTGCCCTGAATGAAGTGCAG
 GACGAGGCAGCGCGGCTATCGTGGCTGGCCACGACGGGCGTTCTTTGCGCAGCT
 GTGCTCGACGTTGTCACTGAAGCGGG
 AAGGGACTGGCTGTCTATTGGGCGAAGTGGCGGGCAGGATCTCCTGTCATCTCAC
 CTGCTCCTGCGAGAAAGTATCCA
 TCTAGGCTGATGCAATGCGGCGGCTGCATACGCTTGATCCGGCTACCTGCCCCATT
 CGACCACCAAGCGAAACATCGCATC
 GAGCGAGCAGCTACTCGGATGGAAGCCGGTCTGTGATCAGGATGATCTGGAC
 GAAGAGCATCAGGGGCTCGCCAGC
 CGAACTGTTCGCCAGGCTCAAGGCGCGCATGCCGACGGCGAGGATCTCGTCTG
 GACCCATGGCGATGCCCTGCTTGCCGA
 ATATCATGGTGGAAAAATGGCCGCTTTCTGGATTATCGACTGTGGCCGGCTGGG
 TGTGGCGGACCGCTATCAGGACATA
 GCGTTGGCTACCCGTGATATTGCTGAAGAGCTTGGCGGCGAATGGGCTGACCGCT
 TCCTCGTGCTTTACGGTATCGCCG
 TCCCGATTTCGACGCGCATCGCCTTCTATCGCCTTCTTGACGAGTTCTTCTGAGCG
 GGAATCTGGGGTTTCGAAATGACCGA
 CCAAGCGACGCCCAACCTGCCATCAGGATGGCGCGCAATAAAATATCTTTATTTTC
 ATTACATCTGTGTGTGTTTCTTGTGTAAGATCCGCGTATGGTGCACTCTC-

Figure 7B

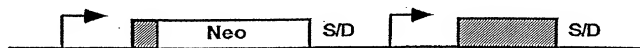
AGTACAATCTGCTCTGATGCCGCATAGTTAAGCCAGCCCCGACACCCGCCAA
 CACCCGCTGACGCGCCCTGACGGGCTTGCTCTGCTCCCGGCATCCGCTTACAGACA
 AGCTGTGACCGTCTCCGGGAGCTGC
 ATGTGTGAGAGGTTTTCACCGTTCATCACCGAAACGCGGAGACGAAAGGGCCTCG
 TGATACGCCTATTTTTATAGGTTAA
 TGTCATGATAATAATGGTTTTCTTAGACGTCAGGTGGCAGTTTTCGGGGAAATGTG
 CGCGGAACCCCTATTTGTTTTATTT
 TCTAAATACATTCAAATATGTATCCGCTCATGAGACAATAACCTGATAAATGCT
 TCAATAATATTGAAAAAGGAAGAGT
 ATGAGTATTCAACATTTTCGGTGTGCGCCCTTATTCCTTTTTTGCGGCATTTTGCT
 TCCGTTTTTGCTCACCCAGAAAC
 GCTGGTGAAGTAAAGATGCTGAAGATCAGTTGGGTGCACGAGTGGGTTACAT
 CGAACTGGATCTCAACAGCGGTAAGA
 TCCCTGAGAGTTTTGCGCCCGAAGAACGTTTTCCAATGATGAGCATTTTAAAGT
 TCTGCTATGTGGCGGTTATTATCC
 CGTATTGACGCGGGCAAGAGCAACTCGGTGCGGCATACACTATTCTCAGAATG
 ACTTGGTTGAGTACTCACAGTCAAC
 AGAAAAGCATCTTACGGATGGCATGACAGTAAGAGAATTATGCAGTGCTGCCATA
 ACCATGAGTGATAACACTGCGGCCA
 ACTTACTTCTGACAACGATCGGAGGACCGAAGGAGCTAACCGCTTTTTTGCAAA
 CTGGGGGATCATGTAACTCGCCTT
 GATCGTTGGGAACCGGAGCTGAATGAAGCCATACAAACGACGAGCGTGACACC
 ACGATGCGCTGTAGCAATGGCAACAAC
 GTTGCGCAAACTATTAACTGGCGAACTACTTACTCTAGCTTCCCGGCAACAATTA
 ATAGACTGGATGGAGCGGATAAAG
 TTGACGAGCACTTCTGCGCTCGGCCCTCCGCTGGCTGTTTATGTCTGATAA
 ATCTGGAGCCGGTGAGCGTGGGTCT
 CGCGGTATCATTGACGACTGGGGCCAGATGGTAAGCCCTCCCGTATCGTAGTTA
 TCTACACGACGGGGAGTCAAGCAAC
 TATGGATGAACGAAATAGACAGATCGCTGAGATAGGTGCTCACTGATTAAGCAT
 TGGTAACTGTGAGCAAGTTTACT
 CATATATACCTTATGATTGATTTAAACTTCACTTTTAAATTTAAAGGATCTAGGTG
 AAGATCCTTTTTGATAATCTCATG
 ACCAAAATCCCTTAACGTGAGTTTTCGTTCCACTGAGCGTCAGACCCCGTAGAAA
 AGATCAAAGGATCTTCTTGAGATCC
 TTTTTTCTGCGGTAACTGCTGCTTGCAAAACAAAAAACCCGCTACCGAGG
 GTGGTTTTGTTGCGGATCAAGAGC
 TACCAACTCTTTTCCGAAGGTAAGTGGCTTACGAGAGCGCAGATACCAAATAC
 TGTCTTCTAGTGTAGCCGTAGTTA
 GGCCACCATTCAAGAACTCTGTAGCACCGCTACATACCTCGCTCTGCTAATCC
 TGTACCAAGTGGCTGCTGCCAGTGG
 CGATAAGTCTGTCTTACCGGGTTGGAAGTCAAGACGATAGTTACCGGATAAGGCG
 CAGCGGTCCGGCTGAACGGGGGTT
 CGTGACACAGCCGAGCTTGGAGCGAACGACCTACACCGAACTGAGATACTAC
 AGCGTGAGCTATGAGAAAGCGCCACGCTTCCGAAGGGAGAAAGCGGACAGGT
 ATCCGGTAAGCGGCAGGGTGGAAACAGGAGAGCGCACGAGGGAGCTTCCAGG
 GGGAAACGCGTGGTATCTTATAGTCTGTGCGGTTTTGCCACCTCTGACTTGAG
 CGTCAATTTTTGTGATGCTCGTCAG
 GGGGGCGGAGCCTATGGAACACGCGAGCAACGCGGCTTTTTACGGTTCTCTGG
 CCTTTTGTCTGGCCTTTTGTCTACATGGCTCGAC3'

Figure 7C

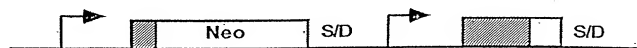
A



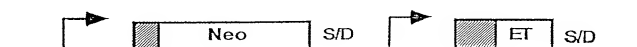
B



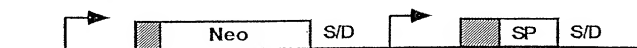
C



D



E



F

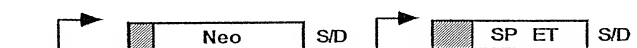


FIGURE 8

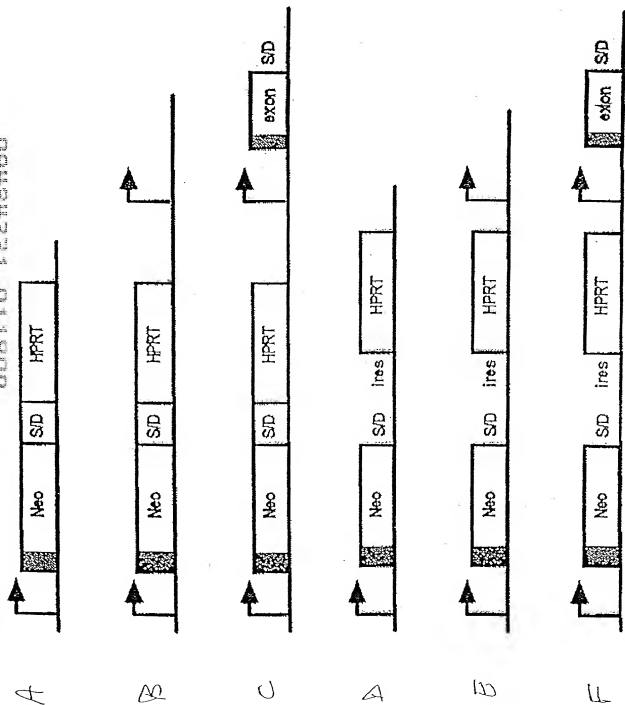


FIGURE 9

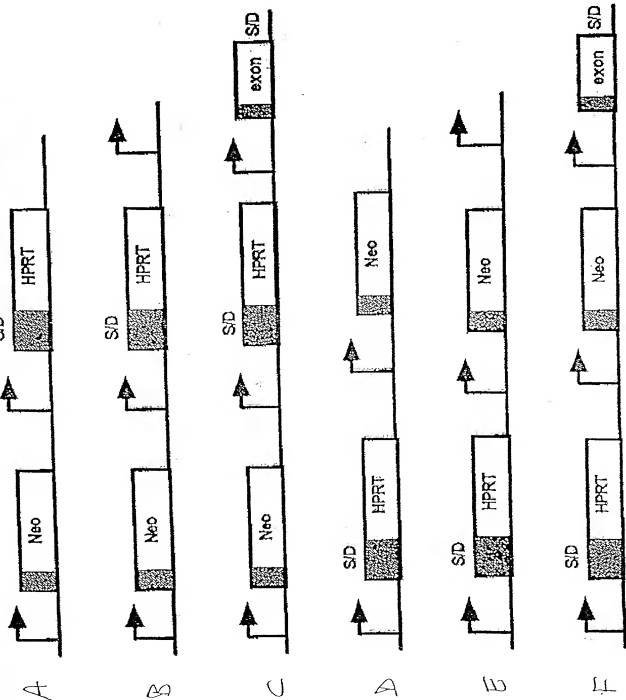


FIGURE 10

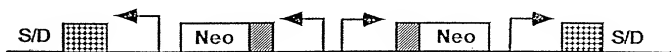
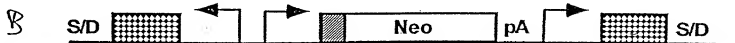
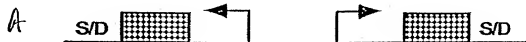


FIGURE 11

008770-1048460

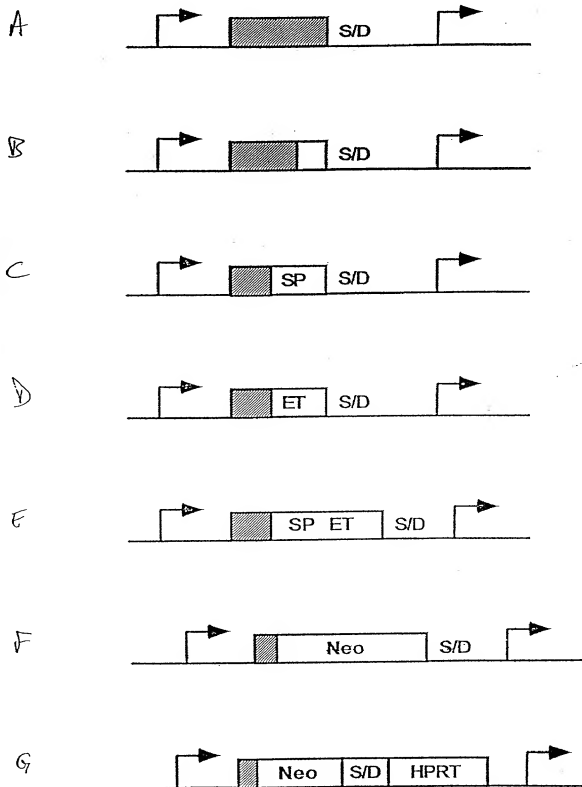


FIGURE 12

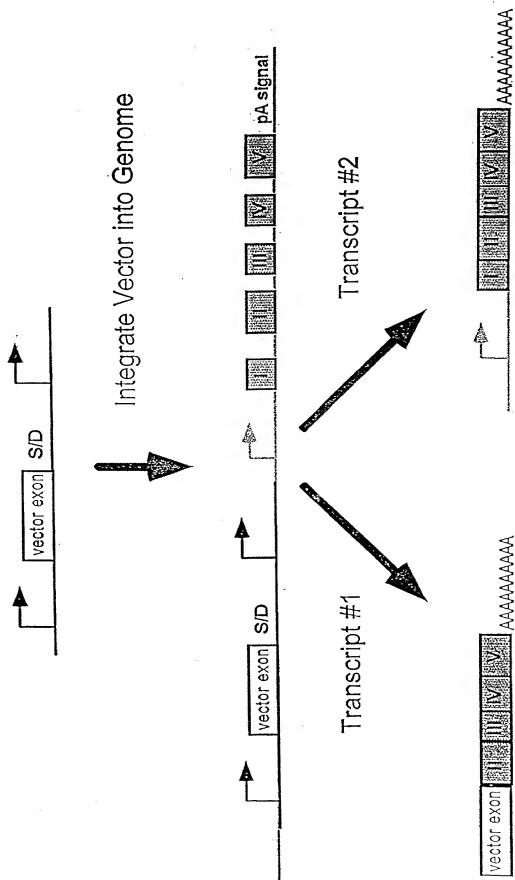


FIGURE 13

AGATCTTCAATATTGGCCATTAGCCATATTATTTCATTGGTTATATAGCATAAATCAATATTGG
 CTATTGGCCATTGCATACGTTGATCTATATCATATATGTACAGTTATATTGGCTCATGTCCA
 ATATAGCCGCCATTGTTGGCATTGATTAITGACTAGTTATTAATAGTAATCAATTAACGGGGTCA
 TTAGTTTCATGCCCATATATGGAGTTCCCGGTTACATAACTTACGGTAAATGGCCCGCTGGC
 TGACGGCCCAACGACGCCGCCCAATTGACGTCAATAATGACGTATGTTCCCATAGTTAACGCCA
 ATAGGGACCTTTCCATTGACGTCATAGGGTGGAGTATTATACGGTAAACTGGCCACTTGGCAGTA
 CACGAGTGTTATCATATGCCAAGTCGCGCCCTATTGACGTCAATGACGTTAAATGGCCCGCT
 TGGCATTATGCCCATACATGACCTTACGGGACTTTCCTACITGGCAGTACATCTACGTATTA
 GTCATCGCTATTACCATGGTGATGCGGTTTGGCAGTACACCAATGGGCGTGGATAGCGGTTT
 GACTCACGGGGATTTCGAAGTCTCCACCCCAATTGACGTCAATGGGAGTTTGTGTTGGACCAA
 AATCAACGGGACTTTCGAAAATGTGCTAACAACTGGGATGGCCGCGCCCGTTGACGCAAAATG
 GCGGTTAGCGGTGATACGGTGGGAGGTCTATATAAGCAGCAGCTCGTTTGGTAAACCGCTCAGAT
 CACTAGAAGCTTTATTGGCGTATGTTTATCACAGTTAAATTTGCTAACGCGAGTCAGTGCTTCTGA
 CACAACAGTCTCGAACTTAAGCTGCAAGTGA CTCTCTTAAATcaacatggctcagggtgactgcGATCTA
 GCGCTATATCGGTTGATGCAATTTCTATGCGCACCCGTTCTCGGAGCACTGTCCGACCGCTTT
 GGCGCGCGGCCAGTCTCTGCTCTGCTACTTGGAGGCCACTATGCACTACGCGATCATGCGGAC
 ACCAACCCCGCTCTGTGGATCTCTACGCGCGGACGCATCGTGGCGCGGCATCACCGCGGCCACA
 GGTGCGGTTGTCTGGCGCTTATATCGCGGACATCACCGATGGGGGAAGTACGGGCTCGGCACTTC
 GGGCTCATGAGCGCTTTTGTGGCTCTCTTAAAGGTAGCAGATCTTGTGCTAGAGTCGCAAAAT
 CTCTATTTTGACACGCTTATCATCGCAGATCTCTGAGCTTGTATGTTGTCATCTCAGTACAATCT
 GCTCTGCTGCGCATAGITAAAGCCAGTATCTGCTCTGCTGTTGTGTTGGAAGTCTGCTGAGT
 AGTGGCGGACGCAAAATTTAAGCTACACAAGGCAAGGCTTGACCGACAATTGCAATGAAGAAT
 TCTGCTTAGGGTTAGGCTTTTGGCGCTGCTTCGCGATGTACGGGCGAGTATACGCGTATCTGA
 GGGGACTAGGGTGTGTTTAGGCGGCCAGCGGGCTTCGGTTGTA CCGCGTTAGGAGTCCCTCTC
 AGGATATAGTAGTTTCGCTTTTGCAATAGGGAGGGGGAAATGTAGTCTTATGCAATACACTTGT
 AGTCTTGCAACATGGTAACGATGAGTTAGCAACATGCGCTTACAAGGAGAGAAAAAGCAACCGT
 GCATCGCGATTTGGTGGAAGTAAGGTGTTACGATCGTGCCTTATTAGGAAGGCAACAGACAGG
 TCTGACATGGATTGGACGAACCACTGAATTCGCGATTCGACAGATAATTGTATTAAAGTGCTC
 AGCTCGATACAATAAACGCCATTGACCATTACCACATTTGGTGTGCCACTCCAAAGCTGGGTA
 CCAGCTGCTAGCCTCGAGACGCGTATTTCCTTCGAAGCTTtcatggttgggtcgtcaaacctcgtcgtggtc
 ccgaacatgggcatcgcaagaacgggacctgcoctggccacgctcagggaatgaattcagatattccagagaatgaccacaacccctcagtaga
 aggtaaacagaatctgggattatgggtaagaagaactggtctccatctcagtagaagaatcagacatttaagggttagaattatgattctcagcagaga
 ctcaagggaacctccacaaggagctcaatttctccagaaggtcagtagatgctttaaactcagacaacagaattagcaataaagtacagatggtt
 ggtatggttgggaggttcttataaagggaagcgaatcaaccaggccatctaaacttattgtgacaaggaatcagcaagcttggaaagtgaacggtt
 ttccagaagaattgatttgaagaataaactctgccagaataaccagggttctctcgtatgtagcaggaggaagaaggcaataagtacaagaattgaagtata
 tgaagaagaattgattatCGATCTTAAAGTTTAACTCTTTCCCGGGGGTACCGTGACCTGCGGCGCGGAATTC
 CAAAGCTTGAGTATTCTATCGTGTACCTAAATAAATCTTGGCGTAATCATGGTTCATATCTGTTTCC
 TGTGTGAAATTGTTATCCGCTCAAAATTCACACACAACATACGAGCCGGAAGCATATAAGTGA
 AAGCCCTGGGTTGCCTAATGAGTGAGCTAACTCAACATTAATTGCGTTTGGCGCATGCTCTCAAAATTT
 TGTGAGGTTTAATGTTCTCGAAGAAGACATGATAAGATACATTGATGAGTTTGGACAAAAACACA
 ACAAGAATGCAAGTGAATAAATGCTTTATTGTTGAAATTTGTGATGCTATTGCTTTATTGTA
 ACCATTATAAGCTGCAATAAACAAGTTAAACAACAACAATTGCATTCATTATTATTTTCAGGTT
 CAGGGGGAGATGTGGGAGGTTTFTAAAGCAAGTAAACCTCTACAAATGTGTTGTAATAACCG
 ATAAGGATCGATTCCCGAGCCTGAATGGCGAATGGACGCGCCTGTAGCGCGCATTAAGCG
 CGCGGGTGTGGTGGTTACGCGCACGTGACCGCTACACTTGGCGAGCGCTTGTAGTGTAGCGCCGCTCC
 TTTGCTTTCTTCCCTTCTTCTCGCCACGTTTCGCGCGCTTTCCCGCTCAAAGCTCTAAATCGG
 GGGCTCCCTTTAGGGTTCCGATTAGTGTCTTACGGCACCTGACGCCCAAAAAAATCTGATTAG
 GGTGATGGTTACAGTAGTGGCCATCGCCCTGATAGACGGTTTTTCGCCCTTTGACGTTGGAG
 TCCACGTTCTTTAATAGTGGACTCTTGTTCGCAACTGGAACAACACTCAACCCATATCTTCGGTC
 TATTCTTTTGTATTAAGGGATTTTTGGCGATTTCGGCTTATTGTTTAAAAAATGAGCTGAATTT
 AACAAAAATTTAACGCGAATTTTAAACAAAAATTAACGCTTTACAAATTTGCGCTGTGACCTTC
 TGAGGCGGAAAGAACCAAGCTGGGAATGTGTGTCAGTTAGGGGTGGAAAGTGTCCGAGGCTC
 CCCAGCAGGCAGAAAGTATGCAAAGCATGCATCTCAATTAGTCAGCAACACAGGTGTGGAAGT
 CCCAGGCTCCCCAGCAGGCAGAAAGTATGCAAAGCATGCATCTCAATTAGTCAGCAACCATAT-

GTCCGCCCTTAACCTCCGCCATCCGCCCTAACCTCCGCCAGTTCGCCCTTCTCCGCC
 ATGGCTGACTAATTTTTTATTTATTTATGACAGAGGCCGAGGCCCTCGGCTCTGAGCTATTTC
 AGAAGTAGTGTAGGAGGCTTTTTTGGAGGCTTAGGCTTTTGCAAAAGGCTTGATTTCTTCGACA
 CACAGTCTCGAACTTAAAGGCTAGAGCCACCATGATTGAACAAGATGGATTGCAACGAGGTT
 CTCGCCGCCCTTGGGTGGAGAGGCTATTCCGGCTATGACTGGGCACACAGACAATCGGCTGC
 TCTGATGCCGCCGTGTTCGGGCTGTACGCGCAGGGGCGCCCGGTTCTTTTGTCAAGACCGAC
 CTGTCCGGTGGCTTAATGAACCTGCAGGACGAGGCGAGGCTATCGTGGGCTGGCCACCGAC
 GGGCGTTCTTTCGCGAGCTGTGCTCGACGTTGTCACTGAAGCGGGAAGGACGCTGGCTGCTATT
 GGGCGAAGTGGCGGGGCGAGGATCTCTGTGATCTCACCTTGCTCTGCGGAGAAAGTATCCAT
 CATGGCTGATGCAATGCGCGGCTGTCATACGCTTGATCCGGCTACCTGCCCATTCGACCAACA
 AGCGAAGCATGCGATCGAGCGAGCAGTACTCGGATGGAAGCGGCTCTTGTCGATCAGGATG
 ATCTGGAACGAAGAGCATCAGGGGCTCGCGCCAGCGCAACTGTTGCGAGGCTCAAGGGCGGC
 ATGCCCGACGGCGAGGATCTCGTCTGACCCATGGCGATGCTGCTTGCGCAATATCATGGTG
 GAAAATGGCGGCTTTTCTGGATTTCATGCACTGTGGCGGCTGGGTGTGGCGGACCGCTATCAG
 GACATACGCTTGGCTACCGGCTATATGCTGAAGAGCTTGGCGGCGAATGGGCTGACCGCTTC
 CTGCTGCTTTACGGTATCCCGCTCCCGGATTGCGAGCGCATCGCTTCTATCGGCTTCTGACG
 AGTTCCTCTGAGCGGAGCTCTGGGGTTCGAAATGACCGACCAAGCGACGCCCAACTGGCAT
 CAGGATGGCGCAATAAAATATCTTTATTTTATTACATCTGTGTGTGGTGTCTTGTGGGAAG
 ATCCCGGCTATGGTGCACTCTCAGTACAATCTGCTCTGATGGCGGATGTTAAAGCGAGGCCGAG
 CACCGCGCAACACCGGCTGACGCGCCCTGACGCGGCTTGTCTGCTCCCGGCTACCGCTTACAGA
 CAAGCTGTGACCGCTTCCGGGAGCTGTCATGTGTGAGAGGTTTTCACCGTCTACCGCAAGCG
 GCGAGACGAAAGGGGCTCGTGATACGCTTATTTTATAGGTTAATGTGATGATAATAATGGTT
 TCTTAGACGTACGGTGGCCTTTTTCGGGGAATGTGGCGGAAACCCCTTGTGTTATTTTCT
 AAATACATTCAAAATATGTATCCGCTCATGAGACAATAACCGTGATAAATGCTCAAAATATTT
 GAAAAGGAAGAGATATGATATTCAACATTTCCGTGTGCGCTTATTCCTCTTTTTCGCGCAT
 TTTGCTTCTGTTTTTGTCTACCCAGAAACGCTGGTGAAGTAAAGATGCTGAAGATCAGT
 TGGGTGCAACGATGGGTTACATCGAACTGGATCTCAACAGCGGTAAGATGCTTGAGAGTTTTT
 GCCCGGAAGAACGTTTTTCCAATGATGAGCACTTTTAAAGTTCTGCTATGTGGCGGATTAAT
 CCGCTATTGACGCGCGGCAAGAGCAACTCGGTGCGCGCATACACTATTCTCAGAATGACTTGG
 TTGAGTACTACCGAGTACAGAAAGCATCTTACGGATGGCATGACAGTAAGAGAATTATGC
 AGTGTCTGCCATAACCATGAGTGATAACAACCTGCGGCAACTTACTCTGACAACGATCGGAGG
 ACCGAAGGAGCTAACCGCTTTTTTGCACAACATGGGGGATCATGTAACCTCGCTTGTATCGTTG
 GGAACCGGAGCTGAATGAAGCCATACCAAACGACGAGCGTGACACCAAGCTGCTGTAGCAA
 TGGCAACAACGTTGCGCAAACTATTAACCTGGCGCAACTACTTACTTGTCTCGCGCAACAAT
 TAATAGACTGGATGGAGCGGATAAAAGTTGACGAGCACTTCTGAGCTCGGCGCTTCCGGCT
 GGCTGGTTTATTGTGTATAAATCTGGAGCGGCTGAGCGTGGGTCTCGCGGTATCATTTGCAGCA
 CTGGGCGCAGATGGTAAGCCCTCCCGTATCGTAGTTATCTACACGACGGGAGTCAGGCAAC
 TATGGATGAACGAAATAGACAGATCGCTGAGATAGGTGCTCACTGATTAAAGCATTTGTAAC
 TGTCAAGCAAGTTTACTCATATATACTTTAGATTGATTAAAACTTCATTTTAAATTTAAAAAG
 GATCTAGTGAAGATCCCTTTTATAATCTCATGACCAAAATCCCTTAAACGTGAGTTTTCGTT
 CCATGAGCGGTGACAGACCCCTAGAAAAAGATCAAAAGGATCTTCTTGAGATCCCTTTTTTCTGCG
 CGTAATCTGCTGCTTGTGCAAAAAAACCACCGCTACCAAGCGGTGTTGTTTGGCGGATCA
 AGAGCTACCAACTCTTTTCCGAAGGTAACTGGCTTCAGCAGAGCGCATACCAAAATACCTGT
 CCTTCTAGTGTAGCGTAGTTAGGCCACACCTTCAAGAACTCTGTAGCACCGGCTATACATCT
 CGCTCTGCTAATCTGTTTACCACTGGCTGCTGCGAGTGGCGATAAGCTGCTTCTTACCGGTT
 GGACTCAAGACGATAGTTACCGGATAAGGCGCAGCGGTGCGGCTGAAACGGGGGTGCTGTGCA
 CACAGCCGAGCTTGGAGCGAACGACCTTACACCGAATGAGATACACTACAGCGTGAGCTGAT
 GAAAGCGCCACGCTTCCCGAAGGGAGAAAGCGGACAGGATCCGGTAAGCGGACGGGTGCG
 GAACAGGAGAGCGTACAGGGAGCTTCCAGGGGGAACGCTGATCTTTATAGTCTCTGCT
 GGGTTTGTGCCACTCTGACTTGAAGCTGCGATTTTGTGATGCTGCTCAGGGGGCGGAGCCTA
 TGGAAAAACGCCAGCAACGCGGCTTTTTACGGTTCTGGCTTTTGTCTGGCTTTTGTCTCAC
 ATGGCTCGAC

FIGURE 14B

[illegible]

FIGURE 15A

CTATTGGGCGAAGTGCCGGGGCAGGATCTCTGTCTCATCTCACCTTGTCTCTGCGGAGAAAAGTA
 TCCATCATGGCTGATGCAATGCGGCGGCTGCATAACGCTTGATCCGGCTACCTGCCCATTOGAC
 CACCAAGCGAAAACATCGCATCGAGCGAGCACGTA CTGGATGGAAGCCGGTCTTGTCTGATCA
 GGTGATCTGGACGAAGAGCATCAGGGGCTCGCGCCAGCCGAACGTTCGCGCAGGCTCAAGG
 CGCGCATGCCCGACGCGGAGGATCTGTCGTGACCCATGGCGATGCTGCTTGCCTGAATATCA
 TGGTGGAAAAATGGCCGCTTTTCTGGATTTCATCGACTGTGGCCGGCTGGGTGTGGCGGACCGCT
 ATCAGGACATAGCGTTGGCTACCCGTGATATTGCTGAAGAGCTTGGCGGCGAATGGGCTGAC
 CGCTTCCTCGTGCTTTACGGTATCGCCGCTCCCGATTGCGAGCGCATCGCCTTCTATCGCCTTC
 TTGACGAGGcaTTCTgatggagtagCGGCCGCTAACCTGGTTGCTGACTAATTGAGATGCATGCTTT
 GCATACTTCTGCTGCTGGGGAGCCTGGGGACTTTCCACACCTAACTGACACACATTCCACA
 GCTGGTTCTTTCCGCCTCAGAAGGTACACAGGCGAAATTGTAAGCGTTAATATTTTGTAAAA
 TTCGCGTTAAATTTTGTAAATCAGCTCATTTTTTAAACCAATAGGCCGAAATCGGCCAAAAATC
 CCTTATAAATCAAAAAGATAGACCGAGATAGGGTTGAGTGTGTTCAGTTTGGAAACAAGAG
 TCCACTATTAAGAACGTGGACTCCAACGTCAAAGGGCGAAAAACCGTCTATCAGGGCGATG
 GCCCAC

FIGURE 15B

[illegible]

FIGURE 16A

GTATCCATCATGGCTGATGCAATGCGGCGGCTGCATACGCITGATCGGGCTACCTGCCATTC
 GACCACCAAGCGAAACATCGCATCGAGCGAGCACGTAICTGGGATGGAAGCCGGTCTTGTGGA
 TCAGGATGATCTGGACGAAAGAGCATCAGGGGCTCGCGCCAGCCGAACCTGTTGCCAGGCTCA
 AGGCGCGCATGCCGACGGGAGGATCTCGTCGTGACCCATGGCGATGCCTGCTTGCCGAAT
 ATCATGGTGGAAAAATGGCCGCTTTCTGGATTTCATCGACTGTGGCCGGCTGGGTGTGGCGGAC
 CGCTATCAGGACATAGCGTTGGCTACCCGTGATATTGCTGAAGAGCTTGGCGGCGAATGGGC
 TGACCGCTTCTCGTGCCTTAACGGTATCGCCGCTCCCGATTCCGACGCGCATCGCCTTCTATCGC
 CTTCTTGACGAGGcaTTCTGctggatggCTacAGGTGcgagoodggcgtcgtgattagtgatgaacagggtatgacctgattta
 ttitgcatacctaalcattatgctgaggatttggaaaggggtttatcctcattggactaatatggacaggactgaactcttgcctcagatgagatgaaggag
 atggaggccatcacattgtagccctctgtgtgctcaaggggggctataaaatcttctgctgacctgctggattacatcaaaagcactgaatagaatagatgata
 gactcattcctatgactgtagattttatcagactgaagagctattgtaatgaccagtcaacaggggacataaaagtattgttggagatgactctcaacttta
 actggaaagaatgtcttgatttggagagataaattgacactggcacaacaaatgcagacttcttctccttggcaggcagataaatacgaatgggcaagg
 tgcgaactgtctgggtaaaaggacccacgaagttgttgatataagccagacttcttggatttgaatitccagacaagtttgtttagatagatgacctga
 ctataatgaatacttcagggtattgaatcatgtttgtgtcattatgtaaaactggaaaagcaaaatacaaaagctaaGCGGCCGCTAACCTGGT
 TGCTGACTAATTGAGATGCATGCTTTGCATACCTTCTGCTGCTGGGGAGCCTGGGGACCTTTC
 ACACCCCTAACTGACACACATTCACAGCTGGTTCCTTCGCTCAGAAAGGTACACAGGCGAAA
 TTGTAAGCGTTAATATTTTGTAAAAATTCGGGTAAAAATTTTGTAAATCAGCTCATTTTTTAA
 CCAATAGGCCGAAAATCGGCAAAATCCCTTATAAATCAAAAGAATAGACCGAGATAGGGPTGA
 GTGTGTTTCCAGTTTGGAAACAAAGAGTCCACTATTAAGAAACGTGGACTCCAACGTCAAAAGG
 CGAAAAACCGTCTATCAGGGCGATGGCCAC

FIGURE 16B

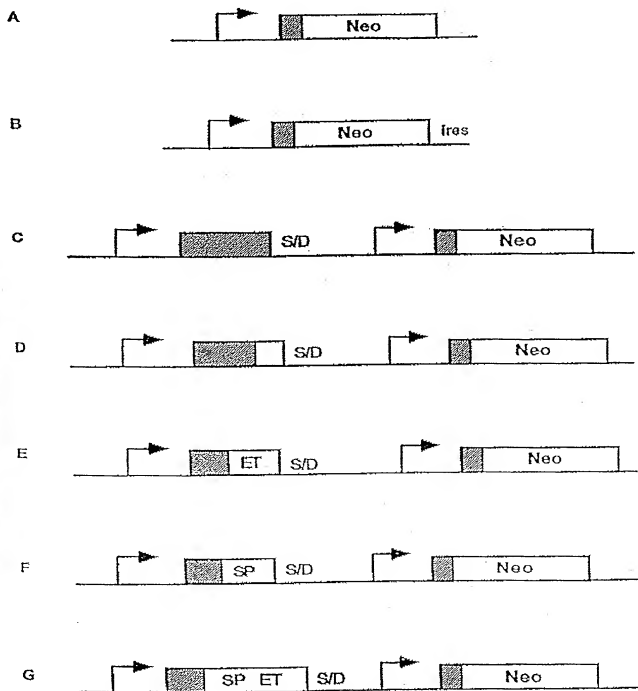


Figure 17

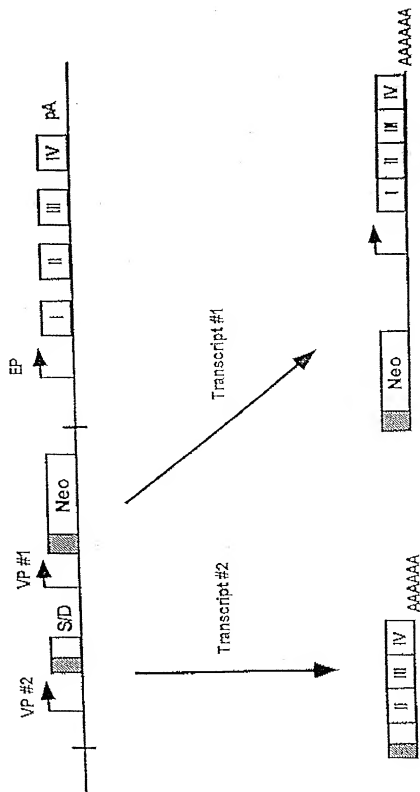


Figure 18

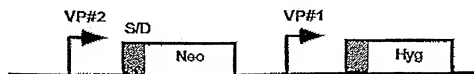


Figure 19

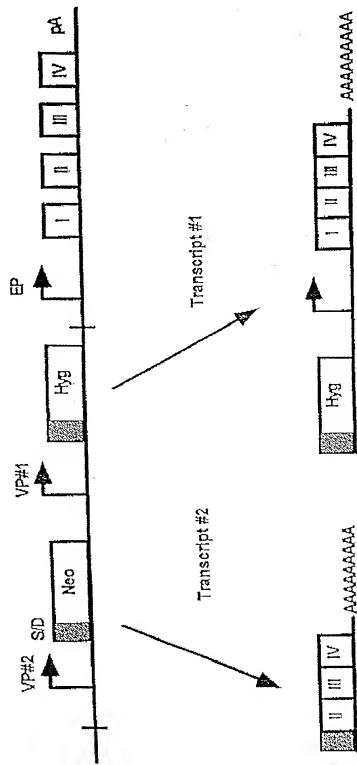


Figure 20A

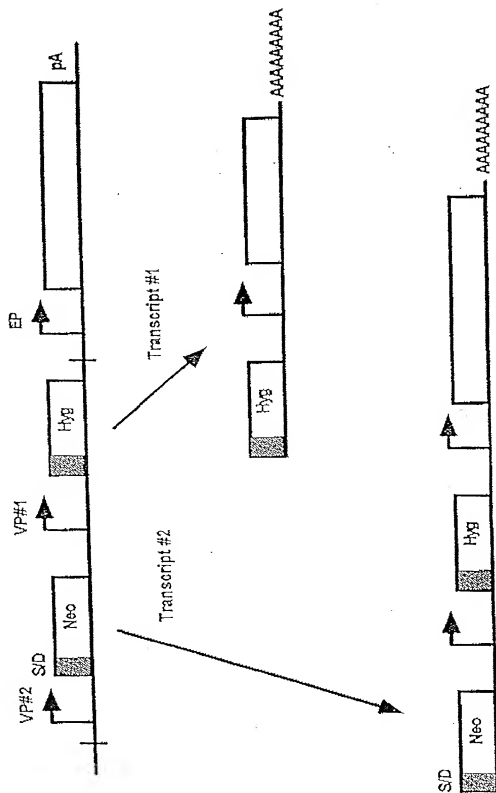


Figure 20B

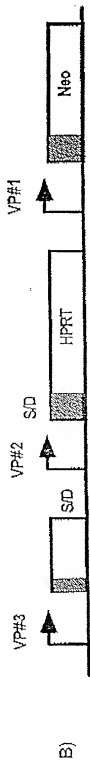
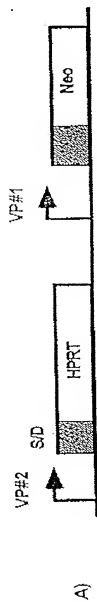


Figure 21

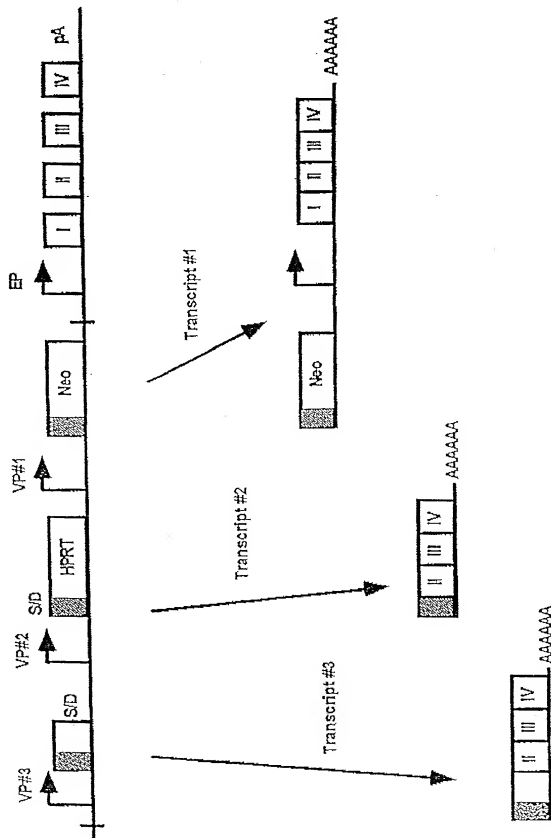
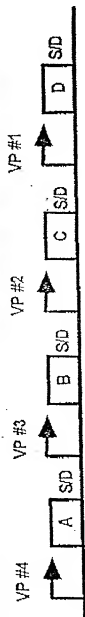


Figure 22



A) Exon A and Flanking Intron

5' UTR	ACCAGGTGATG	Vector Intron
--------	-------------	---------------

B) Exon B and Flanking Intron

5' UTR	ACCATGGCAGGTGATG	Vector Intron
--------	------------------	---------------

C) Exon C and Flanking Intron

5' UTR	ACCATGGCAGGTGATG	Vector Intron
--------	------------------	---------------

D) Exon D and Flanking Intron

5' UTR	ACCATGGCAGGTGATG	Vector Intron
--------	------------------	---------------

Figure 23

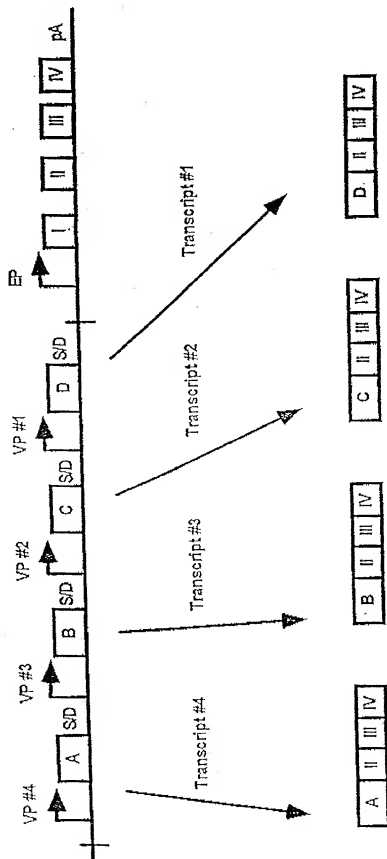


Figure 24

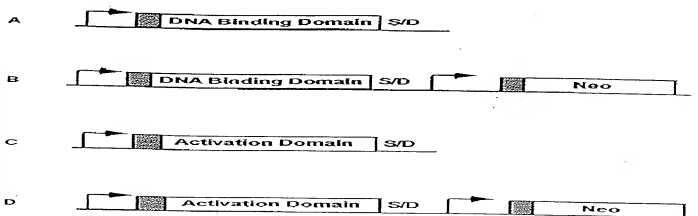


FIGURE 25

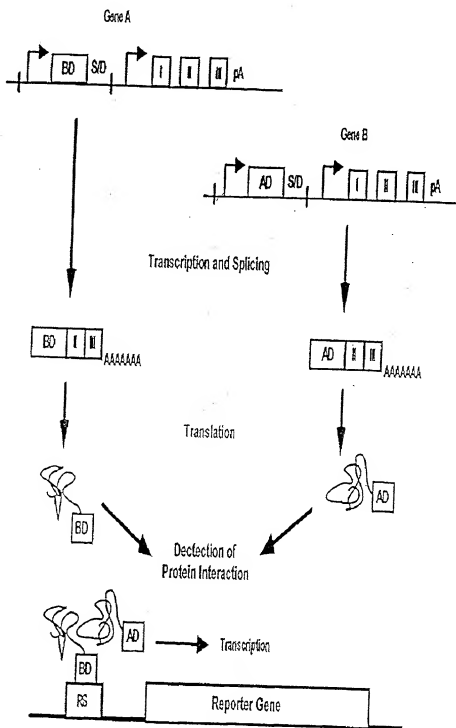


Figure 26

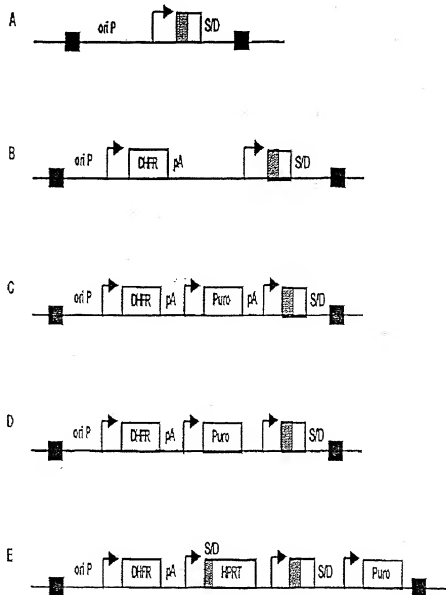


FIGURE 27

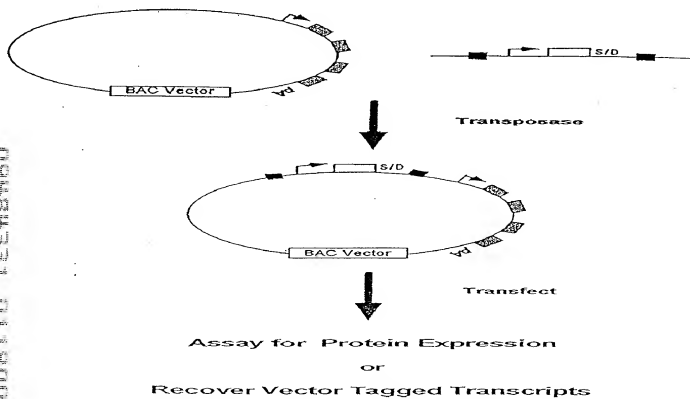


FIGURE 28

CACCTAAATTGTAAGCGTTAATATTTTGTAAAAATTCGCGTTAAATTTTGT
TAAATCAGCTCATTTTTTAACCAATAGGCCGAAATCGGCAAAATCCCTTAT
AAATCAAAAAGAAATAGACCGAGATAGGGTTGAGTGTGTTCAGTTTGGAA
CAAGAGTCCACTATTAAGAACGTGGACTCCAACGTCAAAGGGCGAAAAA
CCGTCTATCAGGGCGATGGCCCACTACGTGAACCATCACCTAATCAAGTT
TTTTGGGGTGCAGGTGCCGTAAAGCACTAAATCGGAACCTAAAGGGAGC
CCCCGATTTAGAGCTTGACGGGGAAGCCGGCGAACGTGGCGAGAAAGGA
AGGGAAGAAAGCGAAAGGAGCGGGCGCTAGGGCGCTGGCAAGTGTAGCG
GTACCGTGCCTGTAAACCAACACACCGCCGCGCTTAATCGCCGCTACAG
GGCGCTCCCATTCGCCATTAGGCTGCGCAACTGTTGGGAAGGGCGATC
GGTGCGGGCCTCTTCGCTATTACGCCAGCTGGCGAAAAGGGGGATGTGCTG
CAAGGCGATTAAAGTTGGGTAACGCCAGGGTTTTCCAGCTCACGACGTTGTA
AAACGACGGCCAGTGAATTGTAATACGACTCACTATAGGGCGAATTGGGT
ACaattcaattcgtcgacctgaaattctaccggtagggaggcgcttttccaaggcagctctggagcatgcgcttag
cagcccgctggggcacttggcgctacacaagtggcctctggcctcgacacattccacatccacggtagggcgcaacc
ggctcgtgttttggggcccttgcgcacacttctactcctccctagtcaggaaagtcccccccgcccgcanctcgcg
tcgtgcaggacgtgacaaatggaaatagcacgtctcactagtctcgtgcagatggacaagcaccgctgagcaatggagc
ggtagggcctttggggcagcggccaatagcagcttggctccttcgcttctgggctcagaggctggnaagggtgggtcc
ggggcgggctcagggcggggctcagggcgggcgggcgccgaaggctcctccggagcggcgccgctcagc
cttcaaaagcgacgtctcggcgctgttctccttctcctcatctccgggcttccgactgcacatcatagatctcagga
gctgaagcttaccatgaccgagtaacaagccacgggtggcctcgccaccccgcgacgacgtccccggcggtacgac
cctcgcccgccgcttcggcgactaccccgccacgcgcacacacgtcgacccggaccgccaatcgagcggggtcaccga
gctgcaagaactctctccacgcgcgtcgggctcgacatcggaagggtgggtcggcgacgacggcgccgctggc
ggttcggacacgcggcgagcgtgcgaagcggggcggtgttcgcccagatcgcccgcgatggcgaggttgagc
gttcccgctggcgccgcgacacacagatggaaaggctcctggcgccgacccggcccaaggcccgctgggtctctt
ggccccacgtcggcgcttctcccgaccacaggcgcaagggtctggcaagcggcctgtgctctccggagtgaggg
cgcgcgagcgcgccggggtgccgccttctggagacctcgccgccccgcaacctcccccttctacgagcggtcgcgctt
caccgtcaccgacgtcgaggtgcccgaaggaccgcgacctgggtgcatgacccgcaagcccggtgcctgacgcc
cgccccagaccgcgagcggccgaccgaaggagcgcacgaccccatgcatcgatggcaggtgaagtatca
aggttagcCATCTTCAATATTGGCCATTAGCCATATTATTCATTGGTTATATAGC
ATAATCAATATTGGCTATTGGCCATTGCATACGTTGTATCTATATCAAT
ATGTACATTTATATTGGCTCATGTCCAATATGACCGCCATGTTGGCATTGA
TTATTGACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTCTCATAGC
CCATATATGGAGTTCGCGTTACATAACTACGGTAATGAGCCCGCTGGC
TGACCGCCCAACGACCCCGCCCATTTGACGTCAATAATGACGTATGTTCCC
ATAGTAACGCCAATAGGGACTTTCATTGACGTCAATGGGTGAGATATTA
CGGTAACACTGCCACTTGGCAGTACATCAAGTGTATCATATGCCAAGTCCG
CCCCCTATTGACGTCATGACGGTAATGGCCCGCTGGCATATAGCCAG
TACATGACCTTACGGGACTTTCCTACTTGGCAGTACATCTACGTATTAGTC
ATCGCTATTACCATGGTGATGCGGTTTTTGGCAGTACACCAATGGGCGTGGGA
TAGCGGTTTGACTACCGGGGATTTCCTAAGTCTCCACCCATTGACGTCAAT
GGGAGTTTGTATTGGCACAAAAATCAACGGGACTTTCAAAAATGTCGTAAAC
AACTCGGATCGCCCGCCCGTTGACGCAAAATGGGCGGTAGGCGGTGTACGG
TGGGAGGCTATATATAAGCAGAGCTGTTTGTAGTGAACCGTACGATCACTAGA
AGCTTTATTGCGGTAGTTTATCAGATTAAATTGCTAACGCAGTCAGTGCT
TCTGACACACAGCTCTGAACTTAAGCTGCAGTGACTCTCTtaataaaccacgcctac
aggtagtactcgATCTGCTACCTTAAGagaggcctatctggcgagtgagcagtgcaagaagaagttaa
GAGAGCCGAAACAAGCGCTCATGAGCCGAAGTGGCGAGCCGATCTTCC
CCATCGGTGATGTCGGCGATATAGGCGCCAGCAACCGCACCTGTGGCGCC-

FIGURE 29A

GGTGATGCCGGCCACGATGCGTCCGGCGTAGAGGATCCACAGGACGGGTG
 TGGTCGCCATGATCGCGTAGTCGATAGTGGCTCCAAGTAGCGAAGCGAGC
 AGGACTGGGCGCGGCCAAAGCGGTCCGACAGTGCTCCGAGAAGCGGGTGC
 GCATAGAAATTGCATCAACGCATATAGCGCTAGATCCTTGCTAGAGTCGAG
 GCCGCCACCGCGGTGGAGCTCCAGCTTTTGTTCCTTTAGTGAGGGTTAAT
 TTCGAGCTTGGCGTAATCATGGTCATAGCTGTTCTGTGTGAAATTGTA
 TCCGCTCACAATTCCACACAACATACGAGCCGGAAGCATAAAGTGTAAG
 CTTGGGGTGCCATAATGAGTGAGCTAACTCACATTAATTGCGTGTCCGCTCAG
 TGCCCGCTTTCCAGTCCGGAAACCTGTCGTGCCAGCTGCATTAATGAATCG
 GCCAACGCGCGGGGAGAGGCGGTTTGGCTATTGGGCGCTCTTCCGCTTCT
 CGCTCACTGACTCGCTGCGCTCGGTCTTCGGCTGCGGCGAGCGGTATCAG
 CTCACTCAAAGGCGTAATACGGTTATCCACAGAATCAGGGGATAACGCA
 GGAAGAACAATGTGAGCAAAAAGGCCAGCAAAAAGGCCAGGAACCGTAAAA
 AGGCCGCGTGTGCTGGCGTTTTTCCATAGGCTCCGCCCCCTGACGAGCATC
 ACAAAAATCGACGCTCAAGTCAGAGGTGGCGAAACCCGACAGGACTATAA
 AGATACCGAGCGTTTCCCCCTGGAAGCTCCCTCGTGCGCTCTCTGTTCGC
 ACCCTGCCGCTTAACGGATACCTGTCGCGCTTTCTCCCTTCGGGAAGCGTG
 GCGCTTTCTCATAGCTCACGCTGTAGGTATCTCAGTTCGGTGTAGGTGCTT
 CGCTCAAAGCTGGGCTGTGTGCAACGAACCCCGTTTCAAGCCGACCGCTGC
 GCCTTATCCGGTAACATATCGTCTTGAGTCCAACCCGGTAAGACACGACTTA
 TCGCCACTGGCAGCAGCCACTGGTAACAGGATTAGCAGAGCGAGGTATGT
 AGCGGCTGCTACAGAGTTCTTGAAGTGGTGGCCTAACTACGGCTACACTAG
 AAGGACAGTATTTGGTATCTGCGCTCTGCTGAAGCCAGTTACCTTCGGAAA
 AAGAGTGGTAGCTCTTGATCCGGCAACAACACCGCTGGTAGCGGTG
 GTTTTTTGTTTGCAAGCAGCAGATTACGGCGAGAAAAAAGGATCTCAAG
 AAGATCCTTTGATCTTTTCTACGGGGTCTGACGCTCAGTGAACGAAAACT
 CACGTTAAGGGGATTTTGGTCAATGAGATTATCAAAAAGGATCTTACCTAGA
 TCCTTTTAAATTAATAATGAAGTTTTAAATCAATCAAAAGTATATGAGT
 AAACCTGGTCTGACAGTTACCAATGCTTAATCAGTGAGGCACCTATCTCAG
 CGATCTGTCTAATTCGTTATCCATAGTTGCCTGACTCCCGCTCGTGTAGAT
 AACTACGATACGGGAGGGCTTACCATCTGGCCCGAGTGCTGCAATGATACC
 GCGAGACCCACGCTCACCGGCTCCAGATTTATCAGCAATAAACAGCCAGC
 CGGAAGGGCCGAGCGCAGAAAGTGGTCTGCAACTTTATCCGCTCCATCCA
 GTCTATTAATTTGTTGCCGGGAAGCTAGAGTAAGTAGTTCGCCAGTTAATAG
 TTTGCGCAACGTTGTTGCCATTGCTACAGGCATCGTGGTGTACCGCTGCTC
 GTTTGGTATGGCTTCAATCAGCTCCGTTCCCAACGATCAAGCGAGTTAC
 ATGATCCCCCATGTTGTGCAAAAAAGCGGTTAGCTCCTTCGGTCTCCGAT
 CGTTGAGCAAGTAAGTTGGCCGAGTGTTATCACTCATGGTATGGCAGC
 ACTGCATAATTCTTACTGTCTATGCCATCCGTAAGATGCTTTTCTGTGACT
 GGTGAGTACTCAACCAAGTCATTCTGAGAATAGTGTATGCGGCGACCGAG
 TTGCTTGTGCCGCGTCAATACGGGATAAATACCGCGCCACATAGCAGAAC
 TTAAAAAGTGCTCATCATTTGGAACCGTTCTTCGGGGCGAAAACTCTCAAG
 GATCTTACCGCTGTTGAGATCCAGTTCGATGTAACCCACTCGTGACCCCAA
 CTGATCTTACGCATCTTTACTTTTACCAGCGTTTCTGGGTGAGCAAAAAAC
 AGGAAGGCAAAATGCCGCAAAAAAGGGAATAAGGGCGACACGGAATGT
 TGAATACTCATACTCTTCTTTTCAATATTATTGAAGCATTTATCAGGGTT
 ATTGTCTCATGAGCGGATACATATTGAATGTATTAGAAAAATAACAAA
 TAGGGGTTCCGCGCACATTTCCCCGAAAAAGTGC

0943-7226/91/0005-0000\$05.00/0

FIGURE 30A

1. **Introduction**
 2. **Background**
 3. **Methodology**
 4. **Results**
 5. **Discussion**
 6. **Conclusion**
 7. **References**
 8. **Appendix**
 9. **Index**
 10. **Table of Contents**
 11. **Abstract**
 12. **Summary**
 13. **Key Words**
 14. **Keywords**
 15. **Subject Headings**
 16. **MeSH**
 17. **Indexing**
 18. **Classification**
 19. **Numbering**
 20. **Ordering**
 21. **Grouping**
 22. **Labeling**
 23. **Marking**
 24. **Notation**
 25. **Symbolism**
 26. **Diagramming**
 27. **Flowcharting**
 28. **Mapping**
 29. **Charting**
 30. **Graphing**
 31. **Tablemaking**
 32. **Formmaking**
 33. **Diagrammaking**
 34. **Flowchartmaking**
 35. **Mappingmaking**
 36. **Chartmaking**
 37. **Graphmaking**
 38. **Tablemaking**
 39. **Formmaking**
 40. **Diagrammaking**
 41. **Flowchartmaking**
 42. **Mappingmaking**
 43. **Chartmaking**
 44. **Graphmaking**
 45. **Tablemaking**
 46. **Formmaking**
 47. **Diagrammaking**
 48. **Flowchartmaking**
 49. **Mappingmaking**
 50. **Chartmaking**
 51. **Graphmaking**
 52. **Tablemaking**
 53. **Formmaking**
 54. **Diagrammaking**
 55. **Flowchartmaking**
 56. **Mappingmaking**
 57. **Chartmaking**
 58. **Graphmaking**
 59. **Tablemaking**
 60. **Formmaking**
 61. **Diagrammaking**
 62. **Flowchartmaking**
 63. **Mappingmaking**
 64. **Chartmaking**
 65. **Graphmaking**
 66. **Tablemaking**
 67. **Formmaking**
 68. **Diagrammaking**
 69. **Flowchartmaking**
 70. **Mappingmaking**
 71. **Chartmaking**
 72. **Graphmaking**
 73. **Tablemaking**
 74. **Formmaking**
 75. **Diagrammaking**
 76. **Flowchartmaking**
 77. **Mappingmaking**
 78. **Chartmaking**
 79. **Graphmaking**
 80. **Tablemaking**
 81. **Formmaking**
 82. **Diagrammaking**
 83. **Flowchartmaking**
 84. **Mappingmaking**
 85. **Chartmaking**
 86. **Graphmaking**
 87. **Tablemaking**
 88. **Formmaking**
 89. **Diagrammaking**
 90. **Flowchartmaking**
 91. **Mappingmaking**
 92. **Chartmaking**
 93. **Graphmaking**
 94. **Tablemaking**
 95. **Formmaking**
 96. **Diagrammaking**
 97. **Flowchartmaking**
 98. **Mappingmaking**
 99. **Chartmaking**
 100. **Graphmaking**
 101. **Tablemaking**
 102. **Formmaking**
 103. **Diagrammaking**
 104. **Flowchartmaking**
 105. **Mappingmaking**
 106. **Chartmaking**
 107. **Graphmaking**
 108. **Tablemaking**
 109. **Formmaking**
 110. **Diagrammaking**
 111. **Flowchartmaking**
 112. **Mappingmaking**
 113. **Chartmaking**
 114. **Graphmaking**
 115. **Tablemaking**
 116. **Formmaking**
 117. **Diagrammaking**
 118. **Flowchartmaking**
 119. **Mappingmaking**
 120. **Chartmaking**
 121. **Graphmaking**
 122. **Tablemaking**
 123. **Formmaking**
 124. **Diagrammaking**
 125. **Flowchartmaking**
 126. **Mappingmaking**
 127. **Chartmaking**
 128. **Graphmaking**
 129. **Tablemaking**
 130. **Formmaking**
 131. **Diagrammaking**
 132. **Flowchartmaking**
 133. **Mappingmaking**
 134. **Chartmaking**
 135. **Graphmaking**
 136. **Tablemaking**
 137. **Formmaking**
 138. **Diagrammaking**
 139. **Flowchartmaking**
 140. **Mappingmaking**
 141. **Chartmaking**
 142. **Graphmaking**
 143. **Tablemaking**
 144. **Formmaking**
 145. **Diagrammaking**
 146. **Flowchartmaking**
 147. **Mappingmaking**
 148. **Chartmaking**
 149. **Graphmaking**
 150. **Tablemaking**
 151. **Formmaking**
 152. **Diagrammaking**
 153. **Flowchartmaking**
 154. **Mappingmaking**
 155. **Chartmaking**
 156. **Graphmaking**
 157. **Tablemaking**
 158. **Formmaking**
 159. **Diagrammaking**
 160. **Flowchartmaking**
 161. **Mappingmaking**
 162. **Chartmaking**
 163. **Graphmaking**
 164. **Tablemaking**
 165. **Formmaking**
 166. **Diagrammaking**
 167. **Flowchartmaking**
 168. **Mappingmaking**
 169. **Chartmaking**
 170. **Graphmaking**
 171. **Tablemaking**
 172. **Formmaking**
 173. **Diagrammaking**
 174. **Flowchartmaking**
 175. **Mappingmaking**
 176. **Chartmaking**
 177. **Graphmaking**
 178. **Tablemaking**
 179. **Formmaking**
 180. **Diagrammaking**
 181. **Flowchartmaking**
 182. **Mappingmaking**
 183. **Chartmaking**
 184. **Graphmaking**
 185. **Tablemaking**
 186. **Formmaking**
 187. **Diagrammaking**
 188. **Flowchartmaking**
 189. **Mappingmaking**
 190. **Chartmaking**
 191. **Graphmaking**
 192. **Tablemaking**
 193. **Formmaking**
 194. **Diagrammaking**
 195. **Flowchartmaking**
 196. **Mappingmaking**
 197. **Chartmaking**
 198. **Graphmaking**
 199. **Tablemaking**
 200. **Formmaking**
 201. **Diagrammaking**
 202. **Flowchartmaking**
 203. **Mappingmaking**
 204. **Chartmaking**
 205. **Graphmaking**
 206. **Tablemaking**
 207. **Formmaking**
 208. **Diagrammaking**
 209. **Flowchartmaking**
 210. **Mappingmaking**
 211. **Chartmaking**
 212. **Graphmaking**
 213. **Tablemaking**
 214. **Formmaking**
 215. **Diagrammaking**
 216. **Flowchartmaking**
 217. **Mappingmaking**
 218. **Chartmaking**
 219. **Graphmaking**
 220. **Tablemaking**
 221. **Formmaking**
 222. **Diagrammaking**
 223. **Flowchartmaking**
 224. **Mappingmaking**
 225. **Chartmaking**
 226. **Graphmaking**
 227. **Tablemaking**
 228. **Formmaking**
 229. **Diagrammaking**
 230. **Flowchartmaking**
 231. **Mappingmaking**
 232. **Chartmaking**
 233. **Graphmaking**
 234. **Tablemaking**
 235. **Formmaking**
 236. **Diagrammaking**
 237. **Flowchartmaking**
 238. **Mappingmaking**
 239. **Chartmaking**
 240. **Graphmaking**
 241. **Tablemaking**
 242. **Formmaking**
 243. **Diagrammaking**
 244. **Flowchartmaking**
 245. **Mappingmaking**
 246. **Chartmaking**<

FIGURE 30B

TTTTGTGTTGCAAGCAGCAGATTACGCGCAGAAAAAAGGATCTCAAGAA
 GATCCTTTGATCTTTTCTACGGGTCTGACGCTCAGTGGAACGAAAACTCA
 CGTTAAGGGATTTTGGTCATGAGATTATCAAAAAGGATCTTACCTAGATC
 CTTTTATCGGTGTGAAATACCGCACAGATGCGTAAGGAGAAAAATACCGCAT
 CAGGAATTTGTAAGCGTTAATAATTAGAAGAACTCGTCAAGAAGCGGAT
 AGAAGGCGATGCGCTGCGAATCGGGAGCGGCGATACCGTAAAGCACGAGG
 AAGCGGTGACGCCATTTCGCCGCCAAGCTCTTCAGCAATATACGCGGTAGCC
 AACGCTATGTCTGATAGCGGTCCGCCACACCCAGCGCCGCCACAGTCGATG
 AATCCAGAAAAGCGGCCATTTTCCACCATGATATTCGGCAAGCAGGCATCG
 CCATGGGTACACGAGAGATCCTCGCCGTGCGGCATGCTCGCCTTGAGCCTG
 GCGAACAGTTGCGCTGGCGGAGCCCTGATGCTCTCTGTCAGATCATCC
 TGATCGACAAGACCGGCTTCCATCCGAGTACGTGCTCGCTCGATGCGATGT
 TTCGCTTGGTGGTTCGAATGGGCAGGTAGCCGGATCAAGCGTATGCAGCCG
 CCGCATTGCATCAGCCATGATGGATACTTTCTCGGCAGGAGCAAGGTGAG
 ATGACAGGAGATCCTGCCCGGCACTTCGCCAATAGACAGCCAGTCCCTTC
 CCGCTTCAGTGACAACGTGAGCACAGCTGCGCAAGGAACGCCCGTCGTG
 GCCAGCCACGATAGCCGCGCTGCCTCGTCTTCGAGTTCATTACGGGCACCG
 GACAGGTGCTGTTCGACAAAAAGAACCGGCGCCCTGCGCTGACAGCCG
 GAACACGCGGCATCAGAGCAGCCGATTGTCTGTTGTGCCAGTCATAGCC
 GAATAGCCTCTCCACCCAAGCGCCGGAGAACCCTGCGTGCAATCCATCTTG
 TTCAATCATGCGAAACGATCCTCATCTGTCTCTTGATCAGAGCTTGATCC
 CCTGCGCCATCAGATCCTTGGCGGCGAGAAAGCCATCCAGTTTACTTTGCA
 GGGCTTGTCAACCTTACCAGATAAAAGTGCTCATCTTGGAAAAcattcaattcgt
 cgacctgcaattcaccgggtaggaggcgcttttcccaaggcagctggagcatgcgctttagcagcccgctgggc
 acttggcgctaca caagtggcctctggcctcgcacattccacatccacggtaggcgcaaccggctccgttcttttgg
 ggccctctcgcgccacttctactctcccttagtcagggaagttccccccgccccgcanctcgctcgtgagagcgtg
 acaaatggaatatgacagctctcactagctcgtgacagtggaacgacccgctgagcaatggagcggttagcccttggg
 gcagcggccaatagcagcttctcctctcgttctcgggctcagagcgctgnaaggggtgggtcggggcgggcctcag
 gggcgggctcaggggcgggcgggcgccgaaggctcctcggaggcccggaftctgcacgcttcaaaagcgacgt
 ctgcccgcgtgttctctcttctcatctccgggctttcgacctgcatccatctagatctcgagcagctgaagcttaccatga
 ccgagtacaagcccgaggctcgctcgcgaaccgagcagcgtccccggggcggtacgcaccctcgcccgcgcttgc
 ccgactaccgcccaacgcgccacacgctcgaccgggaccgccaatcgagcggtgaccgagctgcaagaactcttct
 cagcgcgctcgggctcgacatcgccgaaggtgtgggtcgccgacggcgccgctgggtgggctctggaccacggc
 gagagctgtcaaggcgggggcggtgtcgcgagatcgcccgcatggccgagtcgggttcccgctcggcgtcgccgc
 gcagcaacagatggaaggcctcctggcgccgaccggggcccaaggagcccgctgggtcttggcccaaccgtggggc
 gtcttccggccgaccaggcgcaagggtgtgcgaagcgccgtcgtgctccccggagtgaggcgccgagcgcgccg
 gggtcgccgcttctcggagcagctccgcgcccgaacctccctctacgagcggtctgggttcacgctcaccgcccac
 gtcgaggtgcccgaaggaccgacacgtgggtgcatgacccgcaagccgggtcgtgacgccccccacgacccgca
 gaccggcagcgaagaaggagcgacgacccccatgcatgagtgacgtggcgagtgatcaaggttagcGGCCGC
 TAACCTGGTTTGCTGACTAATTGAGATGCATGCTTTGCTATCTTGCATCTTGCCTGCT
 GGGGAGCCTGGGGACTTTCCACACCCTAACTGACACACATTCCACAGCTGG
 TTCTTTCCGGCTCAGAAGGTACACAGGCGAAATTTGAAGCGTTAATATTTT
 GTTAAATATTCGGGTAAATTTTGTAAATCAGCTCATTTTTTAAACCAATAG
 GCCGAAATCGGCAAAATCCCTTATAAATCAAAAGAATAGACCGAGATAGG
 GTTGAGTGTGTTCAGTTTGGAAACAAGAGTCCACTATTAAAGAACGTGGA
 CTCCAACGTCAAGGGCGAAAAACCGTCTATCAGGGCGATGGCCAC

FIGURE 30C

00404331 011400

agcccgtcctacctgcaatatcaggggactgtgtgacgcttgacgatggagtagatttgccctccctgggtttccacctatg
gtggaaggggctccgcggagggtgatgacggagatgacggagatgaaggagggtgatggagatgaggggtgaggaaag
ggcaggagtgatgttaactgttttagggagacgcctcaatcgattaaagccgtgtattccccccgactaaagaataaatccc
cagtagacatcatgctgtgtgtgtgtatttctggccatctgtctgtcaccatttctgtcctcccaacatggggcaattggg
catlcccatgtgtgtacgtcactcagctccgcgtcaacacctctcgcgttggaaaacattagcgacatttacctggtgagc
aatcagacatcgagcggcttttagcctggcctccttaatticacctaagaatgggagcaacacagcatgaggaaaaggaca
agcagcgaanaatcagcccccttgggagggtggcggcatatgcaaaaggatagcactccaactctactactggtgtatcatat
gctgcatglatatgcatgaggatagcatatgctacccggatagcattaggatagcatatactaccagatatagattaggat
agcatatgctacccagatatagattaggatagcctatgctacccagatataaattaggatagcatatactaccagatataga
ttaggatagcatatgctacccagatatagattaggatagcctatgctacccagatatagattaggatagcatatgctacccag
atatagattaggatagcatatgctacccagatattgggtatgatatgctacccagatataaattaggatagcatatactacct
aatctctattaggatagcatatgctacccggatagcattaggatagcatatactaccagatatagattaggatagcatatactacct
ctaccagatatagattaggatagcctatgctacccagatataaattaggatagcatatactaccagatatagattaggatag
catatgctacccagatatagattaggatagcctatgctacccagatatagattaggatagcatatgctacccagatattgg
gtatgatattgctacccatggcaacattagcccaccgtgtctcagcgacotcgtgaaattaggaccaccaacacctgtgtctt
ggcgctcaggcgcgaagtgtgtglaattgtctccagatcgagcaatcgccccctattctggcccggccacctcttattg
caggtattccccgggtgcatattgtgtgtttgtggcgaagtgtgtgaccgcagtggttagcgggtgtacaaatcagccaa
gttattaccaccttattttacgttccaaaacccaggcggcggtgtgtggggctgacgcgtgccccactccacaatttcaaa
aaaaagagtggccactgtctgtttatgggcccctattggcgtggagccccgtttaaatttccggggtgttagagacaacca
gtggagatccgcgtgctgtcggcgctccactctcttccccctgttacaataagagtgtaacaaacatgggttaccctgtcttggccc
tgccgtggacacatcttaataacccagtatcatattgctactaggattatgtgtgtcccatagccataaattcgtgtgatagtg
acatcagcttcttaccgctgttcccccacccatggatttctattgtaaagataattcagaatgtttctatctcactacatgatttatt
ggccaaagggtgtgtgaggggttatattgtgtgtcatagcacaatgcccaactgaacccccctgcacaaattttattctggggg
cgtcacctgaaacctgtttcgagcacctcacatacccttactgttcaacactcagcagttattctattagctaaacgaagg
agaatgaagaagcagcgcaatgattcaggagagttcaactgcgcgcctgtgatcttcagccagctgccttggcagtaaatgt
gttcactaccctcgtgggaatcctgaccccatgtaataaaaacccgtgacagctcatgggggtggagatatacgtgttcttag
gaccttttactaaccttaattcgatagcatatgcttccgttgggttaacatattgctattgaattagggttagcttggatagat
atactactaccgggaagcatatgctacccgtttagggttaacaggggggtctataaacactattgtaatgcccccttggag
ggctccgttatcggtagctacacaggccccctctgattgacgttgggtgtgacctcccgtagtcttctgggccccgtgggaggt
acatgtccccagcatgtgtgtaagagcttcagccaaggttacacataaaggcaattgtgtgtgacgtccacagactgca
aagtctgtctcaggatgaaagccactcagttgtggcaaatgtgcacatccattataaggatgtcaactacagctcagagaac
cccttgtgttgggtcccccccggtgtcatatgtggaacaggggcccagttggcgaagtgtaccaaccaactgaagggtatc
atgcactgccccgaatacaaaaagcgcctcctcgtacaggcgaagaaggcgagatgccccgttagtcagggtttatgtt
cgtccggcgcgggGCCGCCGAAGGCGCGCGGATCCACAGGACGGGTGTGGCT
GCCATGATCGCGTAGTCGATAGTGGCTCCAAGTAGCGAAGCGGAGCAGGAC
TGGGCGCGGCCAAAGCGGTGCGACAGTGTCTCGGAGAAGCGGTGCGCAT
GAAATTGCATCAACGCATATAGCGTAGATCCTTGCTAGAGTCGAGATCTG
TCGAGCCATGTGAGCAAAAGGCCAGCAAAAGGCCAGGAACCGTAAAAAGG
CCGCGTGTGCTGCGGTTTTTCCATAGGCTCCGCCCCCTGACGAGCATACAA
AAAATTGCAGCTCAAGTCAGAGGTGGCGAAACCCGACAGGACTATAAAGA
TACCAGGCGTTTCCCCCTGGAAGCTCCCTCGTGCCTCTCCTGTTCCGACC
CTGCGCCTTACCGGATACTGTCCGCCCTTTCTCCCTTCGGGAAGCGTGGCG
CTTTCTCATAGCTCACGCTGAGGTATCTCAGTTCGGTGTAGGTCGTTGCCT
CCAAGCTGGGCTGTGTGCACGAACCCCGCTTACGCCACCGCTGCGCCT
TATCCGGTAACATATCGTCTTGAGTCCAACCCGTAAGACACGACTTATCGC
CACTGGCAGAGCCACTGGTAACAGGATTAGCAGAGCGAGGTATGTAGGC
GGTGCCTACAGAGTCTTGAAGTGGTGGCCTAACTACGGCTACACTAGAA
GACAGTATTTGGTATCTGCGCTCTGCTGAAGCCAGTTACCTTCGGA AAAAG
AGTTGGTAGCTCTTGATCCGGCAAAACCAACCCGCTGGTAGCGGTGTTT-

F7JME 31B

GATCTTCAATATTGGCCATTAGCCATATTATTCATTGGTTATATAGCATAAA
TCAATATTGGCTATTGGCCATTGCATACGTTGTATCTATCATATAATGTGA
CATTATATTGGCTCATGTCCAATATGACGCCATGTGGCATTGATTATTG
ACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTTCATAGCCCATAT
ATTGGAGTTCCGGCTTTACATAACTTACGGTAATAATGGCCCGCCTGGCTGACCG
CCCAACGACGCCCCGCCCATGTGACGTCAATAATGACGTATGTTCCCATAGTA
ACGCCAATAGGGACTTTCCATTGACGTCAATGGGTGGAGTATTTACGGTAA
ACTGCCCACTTGGCAGTACATCAAGTGTATCATATGCCAAGTCCGCCCCCT
ATTGACGTCAATGACGGTAATAATGGCCCGCCTGGCATTATGCCCATAGTAC
ACCTTACGGGACTTTCTACTTGGCAGTACATCTACGTATTAGTCATCGCT
ATTACCATGGTGATGCGGTTTTTGGCAGTACACCAATGGGCGTGGATAGCG
GTTTGACTACGGGGATTTCGAAGTCTCCACCCCATGACGTCAATGGGGAG
TTTTGTTTTGGCACAATAATCAACGGGACTTTCCAAAATGTCGTAAACAACCTG
CGATCGCCCCGCCCGTTGACGCAAAATGGGCGGTAGGCGGTGACGTGGGA
GGTCTATATAAGCAGAGCTCGTTTGTGAACCGTCAGATCACTGAATCTTG
ACGACCTACTGATTAACGGCCAGATGTTAGCTAGCGCCGCCACCATGGGCC
CTAAAAAGAAGCGTAAAGTCGCCCCCCCGACCGATGTCAGCCTGGGGGAC
GAGCTCCACTTAGACGGCGAGGACGTGGCGATGGCGCATGCCGACGCGCT
AGACGATTTTCGATCTGGACATGTTGGGGGACGGGATTTCCCGGGGCGCG
GATTTACCCCCACGACTCCGCCCCCTACGGCGCTCTGGATATGGCCGACT
TGGAGTTTGAGCAGATGTTTACCGATGCCCTTGAATTTGACGAGTACCGTG
CGGAATTTGACGTGACTCGCTACCTTAAAggctacttggccggttaacagatgtgtataag
agacagctctcttaaGGTAGCCTGTCTCTTATACACATCTagatccttgcctagagtcgaccaattctc
atgtttgacagctctatcatcgcagatcctgagctgtgatgtgtgcactctcagtaacaatctgtctgtcgcgcatgtttaagcc
agatctgtctcctgcttctgtgtgtggaggtcgtgtagtgcgcgagcaaaatttlaagctacaacgaaggcaggtctgac
gcacaatctcatgaagaatctgcttaggggttaggcgtttttgcgcgtctcgcgatgtacggggccagatatacgcgtatctga
ggggactaggggtgtgttttaggcgcgccagcggggccttgcgtgtacgcgcttagggaggtccctcaggatataagttgttgc
ttttcatagggggggggaattgtagtcttatgcaatacactttagtctgtcaacatgtgaacgatgagttagcaacatgcc
ttacaaggagagaaaaagcaccgtgcacgcgatgttgggaagtaaggtgtacgatcgtgccttattaggaaaggcaaca
gacaggtctgacatggttggaagcaaccactgaattccgcatgacagagataattgtatttaagtgcctagctcgatacaata
aaagcccatttgaccattcaccacatttgggtgtgcacctccaagctgggtaccagctgctagcctcgagacgcgtgatttcctt
cgaagctgtgcatgtgtgtgttcgtcaaaactgcacgtcgtcgtgtgtccagaacatgggcacgtgcgaagaacggggacgtgc
cctggccaccctcctagggaatgaattcagatatttccagagaatgaccacaacctcttcagtagaaggtaaacagaatctggt
gattatgggtaagaagacctgggtctcactctcctgagaagaatgcacctttaaagggtgaattatttattgtctcagcagag
aaactcaagggaacctccacaaggagctcattttcttccagaagctcagatgatgccttaaacctactgaacaaccagaatta
gcaataaagtagacatggtctgtagatgtgtggcagttctgtttataagggaagccatgaatcaccaggcccatctaaac
tatttftgacaaggatcatgcaagactttgaaagtgcacagctttttccagaattgatttggagaataataaaccttgcgcag
aataaccaggtgttctctctgacgtccaggaggagaaaggcaattagaagtacaatttgaagtatagagaagaattTAA
TTAAgggcaccaataactgccttaaaaaaattacgccccgcctgcactcatcgcagctactgtgtatctcatttaagcat
ctgcgcagatcgggaagccatcacagacggcatgatgaacctgaatccgacggcgcacagcactgtgcctctgcgtata
atatttgcctcattgtgaaacccggggcggaagaagttgtccatttggccacgtttaaatcaaacctggtgaactcaccag
ggatttggctgagacgaaaaacataattctcaataaacctttagggaataggccaggttttaccgttaacacgccacatctt
gcgaatatatgtgtagaacctgcgggaatcgtcgtgtgtattcactccagagcgaatgaacaggttctgattgtcctatggaa
aacgggtgtacaagggtgcaactatcccatatcaccaggaatcaccgtcttctcatgcccacagggaattccggatgagcattc
atcaggcgggcaagaatgtgaataaaggccggataaaactgtgtctatttttctacggtctttaaaggccgtgaatatcc
agctgaacggtctgtgtataggttaccattgagcaactgactgaatgcctcaaatgttctttagcatgcaattgggtatataca
acgggtgtataccagtgatttttttctccattttagctcttcttagctcgtaaaaatctcgataaactcaaaaatacgcggcgtgag
tgcatttatttcttagtgggaaggttggaaacctctacgtgcgatcaacgtctcattttgcctaaTTAATTAAGG
CGCGCCgctctcctgctagggatcacgttagaanaaggactaccgacgaaggaaacttgggtgcggctgtgtctgat-

Figure 32A

09 **08** **07** **06** **05** **04** **03** **02** **01**

1. *Chlorophyll a* (Chl *a*)
 2. *Chlorophyll b* (Chl *b*)
 3. *Chlorophyll c* (Chl *c*)
 4. *Chlorophyll d* (Chl *d*)
 5. *Chlorophyll e* (Chl *e*)
 6. *Chlorophyll f* (Chl *f*)
 7. *Chlorophyll g* (Chl *g*)
 8. *Chlorophyll h* (Chl *h*)
 9. *Chlorophyll i* (Chl *i*)
 10. *Chlorophyll j* (Chl *j*)
 11. *Chlorophyll k* (Chl *k*)
 12. *Chlorophyll l* (Chl *l*)
 13. *Chlorophyll m* (Chl *m*)
 14. *Chlorophyll n* (Chl *n*)
 15. *Chlorophyll o* (Chl *o*)
 16. *Chlorophyll p* (Chl *p*)
 17. *Chlorophyll q* (Chl *q*)
 18. *Chlorophyll r* (Chl *r*)
 19. *Chlorophyll s* (Chl *s*)
 20. *Chlorophyll t* (Chl *t*)
 21. *Chlorophyll u* (Chl *u*)
 22. *Chlorophyll v* (Chl *v*)
 23. *Chlorophyll w* (Chl *w*)
 24. *Chlorophyll x* (Chl *x*)
 25. *Chlorophyll y* (Chl *y*)
 26. *Chlorophyll z* (Chl *z*)
 27. *Chlorophyll aa* (Chl *aa*)
 28. *Chlorophyll ab* (Chl *ab*)
 29. *Chlorophyll ac* (Chl *ac*)
 30. *Chlorophyll ad* (Chl *ad*)
 31. *Chlorophyll ae* (Chl *ae*)
 32. *Chlorophyll af* (Chl *af*)
 33. *Chlorophyll ag* (Chl *ag*)
 34. *Chlorophyll ah* (Chl *ah*)
 35. *Chlorophyll ai* (Chl *ai*)
 36. *Chlorophyll aj* (Chl *aj*)
 37. *Chlorophyll ak* (Chl *ak*)
 38. *Chlorophyll al* (Chl *al*)
 39. *Chlorophyll am* (Chl *am*)
 40. *Chlorophyll an* (Chl *an*)
 41. *Chlorophyll ao* (Chl *ao*)
 42. *Chlorophyll ap* (Chl *ap*)
 43. *Chlorophyll aq* (Chl *aq*)
 44. *Chlorophyll ar* (Chl *ar*)
 45. *Chlorophyll as* (Chl *as*)
 46. *Chlorophyll at* (Chl *at*)
 47. *Chlorophyll au* (Chl *au*)
 48. *Chlorophyll av* (Chl *av*)
 49. *Chlorophyll aw* (Chl *aw*)
 50. *Chlorophyll ax* (Chl *ax*)
 51. *Chlorophyll ay* (Chl *ay*)
 52. *Chlorophyll az* (Chl *az*)
 53. *Chlorophyll aza* (Chl *aza*)
 54. *Chlorophyll abz* (Chl *abz*)
 55. *Chlorophyll acz* (Chl *acz*)
 56. *Chlorophyll adz* (Chl *adz*)
 57. *Chlorophyll aez* (Chl *aez*)
 58. *Chlorophyll afz* (Chl *afz*)
 59. *Chlorophyll agz* (Chl *agz*)
 60. *Chlorophyll ahz* (Chl *ahz*)
 61. *Chlorophyll aiz* (Chl *aiz*)
 62. *Chlorophyll ajz* (Chl *ajz*)
 63. *Chlorophyll akz* (Chl *akz*)
 64. *Chlorophyll alz* (Chl *alz*)
 65. *Chlorophyll amz* (Chl *amz*)
 66. *Chlorophyll anz* (Chl *anz*)
 67. *Chlorophyll aoz* (Chl *aoz*)
 68. *Chlorophyll apz* (Chl *apz*)
 69. *Chlorophyll aqz* (Chl *aqz*)
 70. *Chlorophyll arz* (Chl *arz*)
 71. *Chlorophyll asz* (Chl *asz*)
 72. *Chlorophyll atz* (Chl *atz*)
 73. *Chlorophyll auz* (Chl *auz*)
 74. *Chlorophyll avz* (Chl *avz*)
 75. *Chlorophyll awz* (Chl *awz*)
 76. *Chlorophyll axz* (Chl *axz*)
 77. *Chlorophyll ayz* (Chl *ayz*)
 78. *Chlorophyll azz* (Chl *azz*)
 79. *Chlorophyll azaa* (Chl *aza*
 80. *Chlorophyll abz* (Chl *abz*)
 81. *Chlorophyll acz* (Chl *acz*)
 82. *Chlorophyll adz* (Chl *adz*)
 83. *Chlorophyll aez* (Chl *aez*)
 84. *Chlorophyll afz* (Chl *afz*)
 85. *Chlorophyll agz* (Chl *agz*)
 86. *Chlorophyll ahz* (Chl *ahz*)
 87. *Chlorophyll aiz* (Chl *aiz*)
 88. *Chlorophyll ajz* (Chl *ajz*)
 89. *Chlorophyll akz* (Chl *akz*)
 90. *Chlorophyll alz* (Chl *alz*)
 91. *Chlorophyll amz* (Chl *amz*)
 92. *Chlorophyll anz* (Chl *anz*)
 93. *Chlorophyll aoz* (Chl *aoz*)
 94. *Chlorophyll apz* (Chl *apz*)
 95. *Chlorophyll aqz* (Chl *aqz*)
 96. *Chlorophyll arz* (Chl *arz*)
 97. *Chlorophyll asz* (Chl *asz*)
 98. *Chlorophyll atz* (Chl *atz*)
 99. *Chlorophyll auz* (Chl *auz*)
 100. *Chlorophyll avz* (Chl *avz*)
 101. *Chlorophyll awz* (Chl *awz*)
 102. *Chlorophyll axz* (Chl *axz*)
 103. *Chlorophyll ayz* (Chl *ayz*)
 104. *Chlorophyll azz* (Chl *azz*)
 105. *Chlorophyll azaa* (Chl *aza*
 106. *Chlorophyll abz* (Chl *abz*)
 107. *Chlorophyll acz* (Chl *acz*)
 108. *Chlorophyll adz* (Chl *adz*)
 109. *Chlorophyll aez* (Chl *aez*)
 110. *Chlorophyll afz* (Chl *afz*)
 111. *Chlorophyll agz* (Chl *agz*)
 112. *Chlorophyll ahz* (Chl *ahz*)
 113. *Chlorophyll aiz* (Chl *aiz*)
 114. *Chlorophyll ajz* (Chl *ajz*)
 115. *Chlorophyll akz* (Chl *akz*)
 116. *Chlorophyll alz* (Chl *alz*)
 117. *Chlorophyll amz* (Chl *amz*)
 118. *Chlorophyll anz* (Chl *anz*)
 119. *Chlorophyll aoz* (Chl *aoz*)
 120. *Chlorophyll apz* (Chl *apz*)
 121. *Chlorophyll aqz* (Chl *aqz*)
 122. *Chlorophyll arz* (Chl *arz*)
 123. *Chlorophyll asz* (Chl *asz*)
 124. *Chlorophyll atz* (Chl *atz*)
 125. *Chlorophyll auz* (Chl *auz*)
 126. *Chlorophyll avz* (Chl *avz*)
 127. *Chlorophyll awz* (Chl *awz*)
 128. *Chlorophyll axz* (Chl *axz*)
 129. *Chlorophyll ayz* (Chl *ayz*)
 130. *Chlorophyll azz* (Chl *azz*)
 131. *Chlorophyll azaa* (Chl *aza*
 132. *Chlorophyll abz* (Chl *abz*)
 133. *Chlor*

FIGURE 32C

GATCTTCAATATTGGCCATTAGCCATATTATTCATTGGTTATATAGCATAAA
 TCAATATTGGCTATTGGCCATTGCATACGTTGTATCTATATCATAATATGTA
 CATTATATTGGCTCATGTCCAATATGACCCCATGTTGGCATTGATTATTG
 ACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTTCATAGCCCATAT
 ATGGAGTTCCCGGTTACATAACTTACGGTAAATGGCCCCGCTGGCTGACCG
 CCCAACGACCCCCGCCCATTTGACGTCAATAATGACGTATGTTCCCATAGTA
 ACGCCAATAGGGACTTTCCTATTGACGTCAATGGGTGGAGTATTACGGTAA
 ACTGCCCACTTGGCAGTACATCAAGTGTATCATATGCCAAGTCGCCGCCCT
 ATTGACGTCAATGACGGTAAATGGCCCCGCTGGCATTATGCCAGTACATG
 ACCTTACGGGACTTTCCTACTTGGCAGTACATCTACGTATTAGTCATCGCT
 ATTACCATTGGTGTACGGGTTTTTGGCAGTACACCAATGGGCGTGGATAGCG
 GTTTGACTCACGGGGATTTCAAAGTCTCCACCCATTGACGCTCAATGGGAG
 TTTGTTTTTGGCACCAAAATCAACGGGACTTTCAAAATGTCGTAACAACTG
 CGATGCCCGCCCGCTTGACGCAAAATGGGCGGTAGGCGGTGACGGTGGGA
 GGTCTATATAAGCAGAGCTCGTTTGTAGTGAACCGTCAGATCGTAATTCG
 ACGACCTACTGATTAACGGCCAGATCTAAGCTAGCTTCTGAAAGATGAAG
 CTACTGTCTTCTATCGAACAAAGCATGCGATATTTGCCGACTTAAAGAGCTC
 AAGTGCTCCAAAGAAAAACCGAAGTGCGCCAAGTGTCTGAAGAACAACTG
 GGAGTGTGCTGCTACTCTCCCAAAACCAAAGGTCTCCGCTGACTAGGGCACA
 TCTGACAGAAGTGAATCAAGGCTAGAAAGACTGGAACAGCTATTCTACT
 GATTTTTCTCGAGAAGACCTTGACATGATTTTGAAGATGGATCTTTACA
 GGATATAAAAGCATTTGTAAACAGGATTATTTGTACAAAGATAATGTGAATAA
 AGATGCCGTACACAGATAGATTGGCTTCAGTGGAGACTGATATGCCTCTAAC
 ATTGAGACAGCATAGAATAAGTGCACATCATCGGAAGAGAGATAGTA
 ACAAAGGTCAAAGACAGCTTGACTGTATCGCCGGAATTCAGGTAGTACTC
 GCTACCTTAAggcctatctggccgtttaaacagatgtgtataagagacagctctcttaaGGTAGCCTGTCT
 TCTTATACACATCTAgatccttgctagtagtcgaccaattctcagtgtagcagctatcgcagatcctgagct
 tglatgtgcactctcagtaacaatctgctctgctgcgcagtagttaaagcagatctgctcctcctgctgtgtgtgtggaggtcgc
 tlgagtgtgocgagcaaaattaaagctacaacaagggaaggtctgacgcgacaattgcatgaagaatctgcttaggggttag
 ggcgtttgctgctgcttcgcatgtacgggcccagatatacgcgctatctgaggggagctagggtgtgttaggcgcgccagcgg
 ggcttcggtgtacgcggttaggagctccctcaggatagtagtlttgcctttgcataggggagggaattgtagtcttatg
 caatacactgttagtctgtcaacatggtaacgatgtagttagcaacatgccttaagaaggagagaaaaagcacgtgcatgcc
 gattgggtggaagtaagggtgtacgatctgctctattagggaaggcaacagacaggtctgacatggttagtgacgaacacct
 gaattccgcattgcagagataattgtattaaagtgcctagctcgatacaataaacgcatttgacacattcaccacattgggtg
 cactcccaagctgggtaccagctgctagcctcgagacgcgtgatttctcagcagagaaactcaagggaacctccacaaggagctcaatt
 atcgtcgtctgtgtcccagaacatgggcatcggcaagaacgggggacclgcccgtggccacgcctcagggaatgaattcagata
 ttccagagaatgaccacaacctctcagtgaagaagtaaacagaatctgggtgatttgggttaagaagacctgttctccattc
 ctgagaagaatcgacactttaaagggttagaattattttagttctcagcagagaaactcaagggaacctccacaaggagctcaatt
 ctttccagaagctcagatgatgccttaaaacttactgaacaaccagaatagcaataaagtagacatggtctggaatgttg
 tggcagttctggttataagggaagccatgaatcacccaggccactttaaacttttggacaaggatcatgcaagacttgaaa
 gatgaacctgaatgcagcaaatgtgatttggagaaataaaacttccgagaataccagggtgttctctctgcatgtccaggagg
 agaaaaggcaattagtaacaatttgaagtatatgagaagaatgTAAATTAAGgggcaccaataactgccttaaaaaaat
 tacgccccgcctgccaactcatgcagtagtctgtgtaaltcattaaagcattctccgcagatgggaagccatcacagacggcat
 gatgaacctgaatgcagccagccagccatcagcaccctgtgcctgctgataataattggccatgttgtaaaaacggggcggaag
 aagttgtccatattggccacgtttaaatcaaaactggtgaaactcaccaggagattggctgtagacgaaaaacataatctcaat
 aaacccctttagggaataaggccaggttttccacgttaacacgccaatctgcgcaatataattgttagaagaactccgggaatctg
 tctgtgtattcaactccagagcagcagtagaagaacgtlttcagttgtctcatggaacacgggtgtaaacacgggtgaacactat
 caccagctcacccgtcttcttattgccatcgggaattccggatgagcatlcatcaggcgggcaagaatgtgaataaaggccgg
 ataaaactgtgtctatttttcttaccggtctttaaaggccgtaataccagctgaacggctgtgtgtataggatcattgagc-

FIGURE 33A

aactgacgaaatgcctcaaaatgtctttacgatgccattgggatatatacaacggtgttatccagtgattttttccattt
 agcttccttagctcctgaaatctcgataactcaaaaaatcacgccggtagtgatctatttcaattgagaaagtggaaac
 tcttagctgcgatcaacgctctcatttttgcctcaaaTTAATTAAGGCGCGCCgctcctctggttaggagtcacg
 tagaaagactaccgacgaaggaacttgggtgcgcggtgtgtgttataggagtagtaagaacctccctttacaaccta
 ggcgaggaaactgccttgcattccacaatgtgtctttacaccattgagtgctcccccttggaaatggccccctggaccgg
 ccacaaacctggcccgctaaggagggtccattgtctgtttttacggctcttttcaaacctcatatttgcctgaggttttgaag
 gatgcatataggacacttgttatgacaagcccgctcctacctgcaatatacagggtgactgtgtgcagcttgacgatggag
 tagatttgcctccctgggttccaccctatgggtgaaggggggtgcgcgagggtgagacggagatgacggagatgaagg
 agtggatggagatgagggtgagggaaggcaggagtgatgtaacttggtaggagacgcctcaatgctattaaagccgtg
 tattccccgcactaaagaataaatccccagtagacatcatgctgtgtgtgtgtatttctggccatctgtctgttccacatt
 tctgtccccaacatggggcaattgggcatacccatgtgttcacgtcactcagctccgcgctcaacacctctcgtgttga
 aaacattagcgacatttacctgtgagcaatcagacatgcgacgccttttagctcgccctccttaaaatcacctaagaatggg
 agcaaccagcatgcagaaaggacaagcagcgaaaattcacgcccccttgggaggtggcgccatagcaaaaggatag
 cactcccatctactactaggttatcatatgctgactgtatatgcatgaggatagcatatgctacccggatagataggata
 gcatatactaccagataataggatagatagatgctacccagatatagattaggatagcctatgctacccagatataat
 aggatagcatatactaccagatataataggatagatagatgctacccagatataataggatagcctatgctacccagat
 atagattaggatagcatatgctacccagatataataggatagatagcatatgctacccagatataataggatagc
 atataaataggatagcatatactaccctaatctctattaggatagcatatgctacccgatacagataggatagcatatact
 acccagatataataggatagcatatgctacccagatataataggatagcctatgctacccagatataataggatagc
 atatactaccagatataataggatagcatatgctacccagatataataggatagcctatgctacccagatataataggatagc
 ggatagcatatgctacccagatataataggatagcatatgctacccagatataataggatagcctatgctacccagatataataggatagc
 aatattaggatagcatatactaccctaatctctattaggatagcatatgctacccgatacagataggatagcctatgctacccagatataataggatagc
 cctatttggcccgcccaacctactatgcaggttattccccgggtgccaattagtggttttggggcaagtgttggaccgag
 gttttagcgggggttacaacagccaagtattacacccctattttacagttccaaaacccgaggcgccggtgtggggcgta
 cgctgtgcctccctcctcaaatcttcaaaaaaagagtgcccaattgtcttgggtatggcccgccattggcggtggagcccggtt
 aattttcgggggtgttagagacaacacagtgagtcggctgctgtcggcgctcactctcttccccctgtttacaataagatgtg
 aacaacatgttccactgtctgtgtccctcctgggacacatctaataaccacagatcatattgctactaggattatgtgtg
 ccatagcataaattcgtgtgagatggacatccagctctttacggctgtccccaccoccatggtattctattgtaagatattc
 agaattgtttcattcctacactagtatttattgccccagggtttgtgagggttattgtgtgtcatagacaatgcccaccactga
 acccccctgccaatttttattctggggggcgctacctgaaacctgttttcgagcactcacatacacttactgttcaacact
 agcagttattctattagctaaacgaaggagatgaagaagcagcggaagatcaggagaggttcactgcccgtcctgtgtc
 ttacagcactgccccgtgactaaaatgggttactacccctgtgtggaatcctgacccatgtaataaaacctgtgacagctcat
 ggggtgggagatagctgttctttaggaccccttttaaacacctaatctgatagcatatgtctccccgttgggaacatagct
 attgaattgggttagtctcccttatatactactaccgggaagcatatgtctaccccttttaggtttacaaggcccgctta
 taacacactattgctaatggcctctgtgagggtcgttactgtgtagctacacagggccccctgtattgacgttgggtgtgaccc
 cgtagtctcctggccccctggagggttacatgtccccagcattgtgttagaagcttcagcaagattgacataaaggc
 aatgttgtgtgtcagttccagactgcaaaagtctgtccaggatgaagccactcagttgttggcaattgtgcacatccattta
 taaggatgtcaactacagtcagagaaaccttgtgtttgggtcccccccggtgtcactgtgtgaacaggggcccgatgttgga
 agttgtacacaaccaactgaagggtattacatgcactgccccgaatacaaaaacaaagcgctcctcgtacagcgaagaag
 ggccagagatgcgtagtcaggttttagttctgcctggcgccggGCGGCCGCAAGGCGCGCCGATCC
 ACAGGACGGGTGTGGTTCGCCATGTATCGCGTAGTTCGATAGTGGCTCCAAGT
 AGCGAAGCGAGCAGGACTGGGCGGCGGCCAAAGCGGTTCGACAGTGCCTC
 GAGAACGGTTCGCATAGAAATTGCATCAACGCATATAGCGCTAGATCCT
 TGCTAGAGTTCGAGATCTGTTCGAGCCATGTGAGCAAAAAGGCCAGCAAAAGG
 CCAGGAACCGTAAAAAGGCCGCGTGTGCTGGCGTTTTTCCATAGGCTCCGCC
 CCCCTGACGAGCATACAAAAATCGACGCTCAAGTCAGAGGTGGCGAAAC
 CCGACAGGACTATAAAGATACCGAGCGTTTCCCGTGAAGCTCCCTCTGTG
 CGCTCTCCTGTTCGACCCCTGCCGCTTACCGGATACCTGTCCGCCTTTCTCC
 CTTCCGGAAGCGTGGCGCTTTCTCATAGCTACCGCTGTAGGTATCTCAGT-

Case	Age	Sex	Duration	Location	Findings	Comments
1	10	M	10 days	Left eye	Small, dark, pigmented lesion	Benign
2	15	F	2 weeks	Right eye	Small, dark, pigmented lesion	Benign
3	20	M	3 weeks	Left eye	Small, dark, pigmented lesion	Benign
4	25	F	4 weeks	Right eye	Small, dark, pigmented lesion	Benign
5	30	M	5 weeks	Left eye	Small, dark, pigmented lesion	Benign
6	35	F	6 weeks	Right eye	Small, dark, pigmented lesion	Benign
7	40	M	7 weeks	Left eye	Small, dark, pigmented lesion	Benign
8	45	F	8 weeks	Right eye	Small, dark, pigmented lesion	Benign
9	50	M	9 weeks	Left eye	Small, dark, pigmented lesion	Benign
10	55	F	10 weeks	Right eye	Small, dark, pigmented lesion	Benign

FIGURE 33C

1. **Introduction**
 2. **Background**
 3. **Method**
 4. **Results**
 5. **Conclusion**
 6. **References**
 7. **Appendix**
 8. **Table 1**
 9. **Table 2**
 10. **Table 3**
 11. **Table 4**
 12. **Table 5**
 13. **Table 6**
 14. **Table 7**
 15. **Table 8**
 16. **Table 9**
 17. **Table 10**
 18. **Table 11**
 19. **Table 12**
 20. **Table 13**
 21. **Table 14**
 22. **Table 15**
 23. **Table 16**
 24. **Table 17**
 25. **Table 18**
 26. **Table 19**
 27. **Table 20**
 28. **Table 21**
 29. **Table 22**
 30. **Table 23**
 31. **Table 24**
 32. **Table 25**
 33. **Table 26**
 34. **Table 27**
 35. **Table 28**
 36. **Table 29**
 37. **Table 30**
 38. **Table 31**
 39. **Table 32**
 40. **Table 33**
 41. **Table 34**
 42. **Table 35**
 43. **Table 36**
 44. **Table 37**
 45. **Table 38**
 46. **Table 39**
 47. **Table 40**
 48. **Table 41**
 49. **Table 42**
 50. **Table 43**
 51. **Table 44**
 52. **Table 45**
 53. **Table 46**
 54. **Table 47**
 55. **Table 48**
 56. **Table 49**
 57. **Table 50**
 58. **Table 51**
 59. **Table 52**
 60. **Table 53**
 61. **Table 54**
 62. **Table 55**
 63. **Table 56**
 64. **Table 57**
 65. **Table 58**
 66. **Table 59**
 67. **Table 60**
 68. **Table 61**
 69. **Table 62**
 70. **Table 63**
 71. **Table 64**
 72. **Table 65**
 73. **Table 66**
 74. **Table 67**
 75. **Table 68**
 76. **Table 69**
 77. **Table 70**
 78. **Table 71**
 79. **Table 72**
 80. **Table 73**
 81. **Table 74**
 82. **Table 75**
 83. **Table 76**
 84. **Table 77**
 85. **Table 78**
 86. **Table 79**
 87. **Table 80**
 88. **Table 81**
 89. **Table 82**
 90. **Table 83**
 91. **Table 84**
 92. **Table 85**
 93. **Table 86**
 94. **Table 87**
 95. **Table 88**
 96. **Table 89**
 97. **Table 90**
 98. **Table 91**
 99. **Table 92**
 100. **Table 93**
 101. **Table 94**
 102. **Table 95**
 103. **Table 96**
 104. **Table 97**
 105. **Table 98**
 106. **Table 99**
 107. **Table 100**
 108. **Table 101**
 109. **Table 102**
 110. **Table 103**
 111. **Table 104**
 112. **Table 105**
 113. **Table 106**
 114. **Table 107**
 115. **Table 108**
 116. **Table 109**
 117. **Table 110**
 118. **Table 111**
 119. **Table 112**
 120. **Table 113**
 121. **Table 114**
 122. **Table 115**
 123. **Table 116**
 124. **Table 117**
 125. **Table 118**
 126. **Table 119**
 127. **Table 120**
 128. **Table 121**
 129. **Table 122**
 130. **Table 123**
 131. **Table 124**
 132. **Table 125**
 133. **Table 126**
 134. **Table 127**
 135. **Table 128**
 136. **Table 129**
 137. **Table 130**
 138. **Table 131**
 139. **Table 132**
 140. **Table 133**
 141. **Table 134**
 142. **Table 135**
 143. **Table 136**
 144. **Table 137**
 145. **Table 138**
 146. **Table 139**
 147. **Table 140**
 148. **Table 141**
 149. **Table 142**
 150. **Table 143**
 151. **Table 144**
 152. **Table 145**
 153. **Table 146**
 154. **Table 147**
 155. **Table 148**
 156. **Table 149**
 157. **Table 150**
 158. **Table 151**
 159. **Table 152**
 160. **Table 153**
 161. **Table 154**
 162. **Table 155**
 163. **Table 156**
 164. **Table 157**
 165. **Table 158**
 166. **Table 159**
 167. **Table 160**
 168. **Table 161**
 169. **Table 162**
 170. **Table 163**
 171. **Table 164**
 172. **Table 165**
 173. **Table 166**
 174. **Table 167**
 175. **Table 168**
 176. **Table 169**
 177. **Table 170**
 178. **Table 171**
 179. **Table 172**
 180. **Table 173**
 181. **Table 174**
 182. **Table 175**
 183. **Table 176**
 184. **Table 177**
 185. **Table 178**
 186. **Table 179**
 187. **Table 180**
 188. **Table 181**
 189. **Table 182**
 190. **Table 183**
 191. **Table 184**
 192. **Table 185**
 193. **Table 186**
 194. **Table 187**
 195. **Table 188**
 196. **Table 189**
 197. **Table 190**
 198. **Table 191**
 199. **Table 192**
 200. **Table 193**
 201. **Table 194**
 202. **Table 195**
 203. **Table 196**
 204. **Table 197**
 205. **Table 198**
 206. **Table 199**
 207. **Table 200**
 208. **Table 201**
 209. **Table 202**
 210. **Table 203**
 211. **Table 204**
 212. **Table 205**
 213. **Table 206**
 214. **Table 207**
 215. **Table 208**
 216. **Table 209**
 217. **Table 210**

FIGURE 33D

tcaacgacaggagcacgacatgctgcacccgtggccaggacccaacgctgcccgagatgcggccgctgcggctgctgg
 agatggccggacgcgatggatattgtctgccaagggttggttgcgcattcacagttctccgcaagaattgattggctccaatt
 ctggagtggtgaattacgtttagcggaggtgcgcgcggtctccattcaggtgcgggtgcgggtctccatgcacgcgcgacg
 caacgcccgggagggacagcaaggatagggcggcgctcataatccatgccaaacccgttccattgtgctgcggagggcgcc
 ataatccgcgctgacgacatcagcgggtccagtgatcgaagttaggctgtaagagccgcgagcgatccttgaagctgtccct
 gatggctgctcatctacgtccctggacagcatggcctgcaacgcgggcatcccgatgcgcgggaagcgagaagaatcat
 aatgggggaaggccatccagcctgcgctgcggaacgccagcaagacgtgagcccgccgctgcggccgcatgcggcgga
 taatggcctgctctccgcggaacgtttggggcggggacagtgacgaaggcttgagcgagggcgctgcaagattccgaat
 accgcgaagcgacaggccagcatcgtgcgctccagcgaaagcggtctcgcgcgcaaaatgacccagagcgtctgcggc
 accgtgctccatcaggttgcatgataaagaagacagtcataaagtgcggcgacgatagctcatgccccgcgccaccgggaagg
 agctgactgggttgaaggctctcaaggccatcggtgcagcgtctcccttatgcgactcctgattaggaagcagcccgatg
 gtagggtgagccggttgagcaccgcgcggcgaaggaaatgggtgcagcgaaggagatggcgcccaacagctccccggcca
 cggggcctgccccataccacgcgcgaacaaggcgtcatgagccgaagtggcgagcccgatctcccccattcgtgat
 gtggcgatataatggccgacacccgcacgtgtggcgccggtgatgcggccacgatgctgcggcgctagaggaatca
 caggacgggtgtgctgcgcgatgctgcgtatgctgctccaagtgcgaagcgagcaggactggggcgggcgcc
 aaagcggtgcgacagtgctccgagaacgggtgcgcatagaaatgcatcaacgcataagcgtacgacgcgcgcatag
 tgactggcgatgctgtgcggtgagcagatcccgcaaggagcccgccagctacggcgataaacaagcctatgcctacag
 catccagggtagctgtgcccaggatgacgagtgacgacattgttagattcatacaggttgctgactgcgttagcaattaa
 ctgtgataaataccgatataagcttaccgatttccacatgatacagccgcgatgtaattgatacaacagctcatgacg
 tcccgggagcagcaagcccgctcagggcgctgcagccgggtgtggcggtgtgcgggctgcttaactatggcgatc
 agagcagattgtactgaggtgcacatatgcgggtgtgaataccgcacagatgcgttaaggagaaataccgcacatcaggc
 gcattgcgcattcaggtgcgcgaactgttgggaaggcgatgcggcgccctctgcgtatccgcagcgatggcgga
 gggcgatgtgctgcgaaggcgaattaaattgggtaacgcggggtttcccgatcagcagctgttgaataaacgacggccagtg
 attcGAGCTCA TACTTCGAATAGGGATAACAGGGTAAATGCGA TggcggcgccaatCG
 CTCCTTAAGTATGACCGctgcTGGCAAAACAGCTATTATGGGTATTATGGGTATG
 GCCCTAGAAAGCTTggcgttaatcatggtcatatgctgtttcgtgtgaaattgttatccgctcacaaatccacac
 aacatacagccgggaagcataaagtgtaaaagcctgggggtgcctaatgagtgcagtaacacataattgtcggtgtgcctca
 ctgcccgtttccagctcgggaacacgtgctgctccagctgcattaatgacccgcgaggtgcgcgcgcgtaaacccccatcc
 gctgaaagtctgcaaaagcctgatgggacataagtccatcagttcaacggaggtctcacagaaggtttttgcgctggatgtg
 gctgcggcgcaaccgggtgcattgttgcgatccggaggtctgatcggttgcgatgctgaacaattatcctgagaataaattg
 ccttggccctttatgaaatgtggaaactgagtggaatgctgttttctgtttaaacagagaagctgctgtttatccctga
 gaagcgaaacgaacagctgggaaaaatccccattatcgtagagatccgattataatctcagagagcctgtgtagcgtttat
 aggaagtgtgtgttctgtcatgctgcctgcaagcgttaacgaaacgattgaaatgcttccaggaaacataaagaattcttg
 tgcggtgtacgtgtgaattgtgagcggatfatgtcagcaatggacagacaacccatgaacagacaacatgatgtgtgtct
 gtctctttacagccagtagtgcctgcggcagtcgagcgacggggcggaagccctcgagtgcgagcgagggaagcaccagga
 acagcacttatataattctgttacacacgatgcttgaaaaaatctcccctgggggttatccacttatccacggggatattttata
 atatttttttttttagattgtattcttcttttttagagcgctttaggccttatccatgctgttctagagaaggtgtgttgacaa
 attgcctctcagtgtagcaaaatcaccctcaaatgacagtcctgtctgtgacaaattgcccataacccctgtgacaaattgccc
 cagaagaagctgttttttcaaaagtattccctgcttattgactctttttatttagtgcacaaatctgaacttcaacttcc
 atggatctgtcatggcggaaacagcggttatcaatcaacaagaacgttaaaatagcccgcgaaatcgtccagtcacaaacgac
 ctacgtgagggcgcatatagctctccgggatacaaaacgtagtgcgtatctgttgaccagatcagaagaatctgatg
 gcacccatacaggaaatcagcggatctgcggagatccatgttgcataatgtcgtgaaatattcggatgactcctgcgggaagc
 cagtaaaagatatcggcaggcattgaagagtttcggcggggaaggaggtttttatgcgcctgaagaggtatgcggcg
 atgaaaaaggctatgaaattcttctgtttatcaaacgtgcgcacagctcatccagagggctttacaggtctgtgttaatgaga
 tcaacagcagaactccaatgcgctctctatacttgagaaaaaagaaggccgcgacgactatctgatttttcttccg
 catatctcaattccctcttctatcggggttacagaaccggtttacgcagtttgcgcttagtgaacaaaaaagaatcacaatccgt
 atgccaatgcgtttatagcaatccctgtgtcagtatcgtgaagccgaggtgcctcagggcatcgtctcttgcataatcgaactggtat
 atagagcgttgcagctgctcaaaagtaccagcgttatgcctgcactccgcgcgcttccctgcaggtctgtgttaatgaga
 tcaacagcagaactccaatgcgctctctatacttgagaaaaaagaaggccgcgacgactatctgatttttcttccg
 cgtatcatcttccatgacgacagagatagctgagggttatctgtcagagattgagggtgtgttcgtacattgtttctgacct-

GATCTTCAATATTGGCCATTAGCCATATTATTCATTGGTTATATAGCATAAA
TCAATATTGGCTATTGGCCATTGCATACGTTGTATCTATATCATATAATATGA
CATTTATATTGGCTCATGTCCAATATGACCGCCATGTTGGCATTGATTATTG
ACTAGTTATTAATAGTAATCAATTACGGGGTCAATTAGTTCATAGCCCATAT
ATGGAGTTCCCGCTTACATAACTTACGGTAAATGGCCCGCTGGCTGACCG
CCCAACGACCCCGCCCAATTGACGTCAATAATGACGTATGTTCCCATAGTA
ACGCCAATAGGGACTTTCATTGACGTCAATGGGTGGAGTATTTACGGTAA
ACTGCCCACTTGGCAGTACATCAAGTGTATCATATGCCAAGTCCGCCCTCT
ATTGACGTCAATGACGGTAAATGGCCCGCTGGCATTATGCCAGTACATG
ACCTTACGGGACTTTCCTACTTGGCAGTACATCTACGTATTAGTCATCGCT
ATTACCATGGTGATGCGGTTTTGGCAGTACACCAATGGCGGTGGATAGCG
GTTTGACTCACGGGATTTCCAAGTCTCCACCCATTGACGTCAATGGGAG
TTTGTTTTGGCACCAAAATCAACGGGACTTTCCAAAATGTCTGAACAACTG
CGATCGCCCGCCCCGTTGACGCCAAATGGGCGGTAGGCGGTGACGGTGGGA
GGTCTATATAAGCAGAGCTGtttagtaacogtcagatcactgaattctgacacogtcactgattaaoggc
catagaggcctcctgcagaactgtcttagtgacaactatCGATTTCCACACATTATACGAGCCGAT
GTTAAITGTCAACAGCTCATGTCATGACGTCCCGGGAGCAGACAAGCCCGacc
atggctcgagTAATACGACTCACTATAGGGCGACAGGTGAGTACTCGTCACTT
AAGAGAGGCCCTATCTGGCCAGTTAGCAGTCCGAAAGAAAGTTTAAAGAGA
GCCGAAACAAGCGCTCATGAGCCGAAGTGGCGAGCCCGATCTTCCCAT
CGGTGATGTCGGCGATATAGGCGCCAGCAACCGCACCTGTGGCGCCGGTG
ATGCCGGCCACGATGCGTCCGCGTAGAGGATCCACAGGACGGGTGTGGT
CGCCATGATCGCGTAGTCGATAGTGGCTCCAAGTAGCGAAGCGAGCAGGA
CTGGGCGCGGCCAAAAGCGGTCCGACAGTGTCTCCGAGAACGGGTGCGCAT
AGAAATTGCATCAACGCATATAGCGTAGATCTTGCTAGATCGAGATCT
GTCGAGCCATGTGAGCAAAAGGCCAGCAAAAGGCCAGGAACCGTAAAAAG
GCCGCGTTGTCTGGCGTTTTTCCATAGGCTCCGCCCCCTGACGAGCATCAC
AAAAATCGACGCTCAAGTCAGAGGTGGCGAAACCCGACAGAGCTATAAAG
ATACCAGGCGTTTTCCCCCTGGAAGCTCCCTCGTGCCTCTCTCTGTTCCGAC
CCTGCCCTTACCGGATACCTGTCCGCTTTCTCCCTTCGGGAAGCGGTGGC
GCTTTCTCATAGCTCACGCTGTAGGTATCTCAGTTCGGTGTAGTCTGTCG
CTCCAAGCTGGGCTGTGTGCACGAACCCCCCGTTACGCCCGACCGCTGCGC
CTTATCCGGTAACCTATCTGCTTGTAGTCCAACCCCGTAAGACACGACTATC
GCCACTGGCAGCCACTGGTAACAGGATTAGCAGAGCGAGGTATGTAG
GCGGTGCTACAGAGTTCTTGAAGTGGTGGCCTAACTACGGCTACACTAGAA
GGACAGTATTGTGATCTGCGCTCTGCTGAAGCCAGTTACCTTCGGA AAAA
GAGTGTGGTAGCTTGTATCCGGCAAAACAAACCCGCTGGTAGCGGTGGTT
TTTTTGTTTGC AAGCAGCAGATTACGCGCAGAAAAAAGGATCTCAAGAA
GATCCTTTGATCTTTTCTACGGGTCTGACGCTCAGTGGAAACGAAAACCTCA
CGTTAAGGGATTTTGGTTCATGAGATTATCAAAAAGGATCTTTCACCTAGATC
CTTTTatcggtgtgaaataccgcacagatgcgtgaaggagaaataccgcacaggaattgtgaagcgtttaataaftcag
aagaactcgtcaagaaggcgatagaaggcgatgcgtcgcaatcgggagcgcgataaccgtaaaagcagaggaaagcg
gtcagcccaattcgccccaagctcttcagcaatcacgggttagccaacgctatgtcctgatagcgttcgccacacccag
ccggccacagctgatgaatccagaaagcggcacatttccaccatgatattcggcaagcaggcatcgccatgggtcacga
cgagatcctcgccgtcgggcatgctcgcttgagccttgccgaacagttcggctggcgagccccgtgatctctctgctcc
agatcatcctgatgcagacagccgcttcacatccagtagtctgctgctcgatgctatgttctcgttggttcgaatgggc
aggtagccgatcaagcgtagtcagccgcccgttcgcatcagccatgatggatactttctcggcaggagcaagggtgagat
gacaggagatcgtccccggcacttcgcccataagcagccagtccttccccgttcagtgaacgctcgagcacagctgc
gcaaggaaagcccgctgtgcccagccagatagccgctgctcgtctgtcagttcaatcagggcacccgacagaggtc-

FIGURE 35A

ggctctgacaaaaagaacggggcgccccctgcgctgacagccgggaacagcgccgcatcagagcagccgattgtctgtt
 gccacgtcatagccgaatagcctctccaccaagcgccggagaaacctggtgcaatccatctgttcaatcatgogaaac
 gatccctcatcctgtctctgacagagcttgatccccctgcgccatcagatccttggcggcgagaaagccatccagtttacttt
 gcaggcgtctgcaacctaccagatAAAAGTGTCTATCATTTGGAAACAGTTCAATTCTGAG
 GCGGAAAGAACACAGCTGTGGAATGTGTGTCACTAGGGTGTGGAAAGTCC
 CCAGGCTCCCCAGCAGGCAGAAAGTATGCAAAAGCATGTCATCTCAATTATGTCA
 GCAACCAGGTTGTGGAAGTCCCCAGGCTCCCCAGCAGGCAGAAAGTATGCA
 AAGCATGCATCTCAATTAGTCAGCAACCATAGTCCCCGCCCTAACTCCGCG
 CATCCGCGCCCCTAACTCCGCCCAGTTCGCCCATTTCTCCGCCCATGGCTG
 ACTAATTTTTTTTATTTATGTCAGAGGCCGAGGCCGCTCGGCCCTGTGAGCT
 ATTCCAGAAGTAGTGAGGAGGCTTTTTTGGAGGCCATAGGCTTTTTGCAAAAAA
 GCTTGATTCTCTTGACACAAACAGTCTCGAACTTAAGGCTAGAGCCACCATG
 ATTGAACAAGATGGATTGCACGCAGGTTCTCCGGCCGCTTGGGTGGAGAG
 GCTATTCCGGCTACTGCTGGGCACAAACAGACAATCGGGCTCGTCTGATGCCGC
 CGTGTTCGGCTGTGTCAGGCAGGGGCGCCGTTCTTTTTGTCAAGACCGA
 CCTGTCCGGTGCCCTGAATGAACTGCAGGACGAGGCAGCGCGGCTATCGT
 GGCTGGCCACAGCGGGCGTTCCTTGCGCAGCTGTGCTCGACGTTGTCACTG
 AAGCGGGAAGGAGTGGCTGTCTATTGGGCGAAGTGGCGGCGAGGATCTC
 CTGTCTCATCTCACCTGCTCTGCGGAGAAAGTATCCATCATGGCTGATGCA
 ATGCGGCGGCTGCATACGCTTGATCCGGCTACCTGCCCATTCGACCACCAA
 GCGAAACATCGCATCGAGCGAGCACGTACTCGGATGGAAGCCGGTCTTGT
 CGATCAGGATGATCTGGACGAAGAGCATCAGGGGCTCGCGCCAGCCGGAAC
 TGTTCCGCCAGGCTCAAGGCGCGCATGCCGACGCGGAGGATCTCGTCTGT
 ACCATGGCGATGCCTGCTTGCCGAATATCATGGTGGAAAAATGGCCGCTTT
 TCTGGATTTCATCGACTGTGGCCGCTGGGTGTGGCCGACCCGCTATCAGGAC
 ATAGCGTTGGCTACCCGTGATATTGCTGAAGAGCTTGGCGGCGAATGGGCT
 GACCGCTTCTCGTGCTTTACGGTATCGCCGCTCCCGATTTCGCAGCGCATC
 GCCTTCTATCGCCTTCTTGACGAGGcaTTCTgtgagcaggttaagtgcgagccgtgagctgtgatt
 agtgaatgatgaaccaggttatgacctgatttatctgacatacctaattcattatgctgaggatttggaaagggtgtttatccctca
 tggactaattatggacagagctgaacgtctgtctgcagatgtgatgaaggagatgggagggccatcacattgatgacctctg
 tgtgtctcaagggggctataaattcttctgacctgtggtgattacataaagcactgaatagaatagtatgatgataccatt
 ctatgactgtgatatttatcagactgaagagctattgtaatgaccagctcaacaggggacataaaagttaattgttgagatgat
 ctctcaactttaactgaaagaatgtcttgattgtggaagataaattgacactggcaaaacaaatgcagacttgccttctgtg
 gtcaggcagatataaagaaatgggtcaaggtgcgaagcttggtgaaaggacccacgaagtgttgatataagcc
 agacttgttgattgaaattccagacaaagtgtgtgtaggataatgacctgactataatgaatacttcagggatttgaatcat
 gtttgtgtcattagtgaactggaaagcaaaatacaaggcctaaGCGGCCGCTAACTGGTGTGCTGA
 CTAATTGAGATGCATGCTTTGCATACTTCTGCTGCTGGGAGGCTGGGA
 CTTTCCACACCCCTAACTGACACACATTTCCACAGCTGGTGTCTTTCCGCGCTCAG
 AAGGTACACAGGCGAAATTGTAAGCGTTAATATTTTGTAAAAATTCCGCGTT
 AAAATTTTGTAAATCAGCTCATTTTTTAAACCAATAGGCCGAAATCGGCAA
 AATCCCTTATAAATCAAAAAGAAATAGACCGAGATAGGGGTTGAGTGTGTGTTCC
 AGTTTGGAACAAGAGTCCACTATTAAAGAACGTGGACTCCAACGTCAAAG
 GGCGAAAAACCGTCTATCAGGGCGATGGCCCCAC

FIGURE 35B

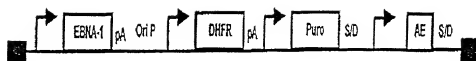


FIGURE 36

GATCTTCAATATTGGCCATTAGCCATATTATTCATTGGTTATATAGCATAAA
 TCAATATTGGCTATTGGCCATTGCATACGTTGTATCTATATCATAATATGTA
 CATTATATTGGCTCATGTCCAATATGACCGCCATGTTGGCATTGATTATTG
 ACTAGTTATTAATAGTAATCAATTACGGGGTCATTAGTTTCATAGCCCATAT
 ATGGAGTTCGCGTTACATAAATTACGGTAAATGGCCGCTGGCTGCGACCG
 CCCAACGACCCCCGCCATTGACGTCAATAATGACGTATGTTCCCATAGTA
 ACGCCAAATAGGACATTTCATTGACGTCAATGGGTGGAGTATTACGGTAA
 ACTGCCCACTTTGGCAGTACATCAAGTGTATCATATGCCAAGTCGCCCCCT
 ATTGACGTCAATGACGGTAAATGGCCGCTGGCATTATGCCAGTACATG
 ACCTTACGGGACTTTTCTACTTGGCAGTACATCTACGTATTAGTCATCGCT
 ATTACCATGGTGATGCGGTTTGGCAGTACACCAATGGGCGTGGATAGCG
 GTTTGACTCACGGGATTTCCAAAGTCTCCACCCCATTTGACGTCAATGGGAG
 TTTGTTTGGCACCAAAATCAACGGGACTTTCCAAATGTCTGTAACAACGTG
 CGATCGCCCGCCCGTTGACGCAAAATGGGCGGTAGGCGTGTACGGTGGGA
 GGTCTATATAAGCAGAGCTCGTTTAAAGTGAACCGTCAATCAATTCG
 ACGACCTACTGATTAAACGGCCATAGAGGCTCTCTGCAGAACTGTCTTAGTG
 ACAACTATCGATTTCACACATTATACGAGCCGATGTTAATTGTCAACAGC
 TCATGCTACGACGTCGCGGAGCAGACAAGCCGACCTAGTCTCGATTAAT
 ACGACTCACTATAGGGCGACAGGTGAGTACTCGCTACCTTAAggcctatctggcgc
 ttttaaacagatgtgtataagagacagctctcttaaGGTAGCTGTCTTTATACACATCTAgatccttg
 cttagctgcaccaattctcatgtttgacagcttatcatcgacagctcgtgactgtgacactcagtaacaatctgctct
 gctgcgcagatagtaagccagatctctccctcgttgtgtgtggaggtcgtgagtagtgcgcgcagcaaaatttaagcta
 caacaggccaggcttgaccgacaattgcatgaagaatctgcttaggggtaggcgttttgcgcgtcttcgcgatgacagg
 ccagatatacgcgtatctgaggggactagggtgtgtttaggcgccagcggggcttcgtgtgacgcgtttaggagtc
 ctacagatagtagtttgcctttgcatagggagggggaaatgtagtcttatgcaatacacctgtgactgtgcaacatctgta
 cgtgaggttagcaacatgccttcaaggagagaaaaagcaccgtgcatgcccgttgggaagtaaggtgtgacgatcgt
 gccctattaggaagcgaacagacaggtctgacatggttggacgaaccactgaattccgcattgcagagataattgtatta
 agtgccttagctcgatacaataaacgccatttgaccattcacacatctgtgtgcacctccaagctgggtaccgctcctagc
 ctgcagacgcgtgatttcttgcgaagcttgcattggtgtgttcgtctaaactgcacgtgcgtgtgtgtcccaaacatggggc
 ggcgaagacggggacgtgcctgcgcacccgtcagggaatgaattcagatattccagagaatgaccacaacctctcagt
 agaaggttaaacagaatctgtgtgatttggtgtaagaagacgtgttctccattcctgagaagaatcgacattaaagggtga
 attaatgtgtctgcagagaaatcgaaggacctccacaaggagctcattttcttcagaagtctagatgatgccttaaaa
 cttaactgaacacagaatttagcaataaagttagacatggtctgtagtgggtggcagttctgtttataagggaagccatga
 atcacccaggcccatcttaaacatttggacaaggatcatgcaagactttgaaagtgcacgttttttccagaattctattgg
 agaaataataaacttgcgcagaataccagggtgttctctctgatgtccaggaggagaaagcgaattagtaacaaattggaat
 atatgagaagaattTAAATTAAGggcaccataaactgccttaaaaaaattacgcccgcctgcacatcgcagat
 actgtgttaattcattgaactctctgcacatggaagccatcacagacgcgatgatgaacatctcaatccgacgcgcatca
 gcacactgtgcctgctgataaatttggccatggtgtaaaacggggcggaagaagtgtccatattggccacgttttaaatca
 aaactggtgaaactcacccagggttggctgagacgaaaaacatatctcaataaaaccttagggaaataggccaggtttt
 caccgtgaacgcgcacactcttgcgaatatatgtgtagaactgcggaaactcgtgtgtgattcaactccagagcgaatgaa
 acgtgtcattgtctcatggaaacgggtgtgaacagggtgaacactatcccatcaccagctcacgtcttcttcatgtccata
 gcgaattccggtgatgacattcaccggcggaagaatgtgaataaaggccgggataaaactgtgtctattttcttaccggt
 cttaaaaggccgctaataaccagctgaacggctcgtttataggtagactagcaactgcgaagtgcctcctaataatgttctt
 acgatgcatttgggataatacaacgggtggtatccagtgatttttttctccattttagcttccatgacctgaataactctgata
 actcaaaaaatacgcgcgtagtgatctatttctattatggtgaagattgggaacctcttaactgcgtcgaacgtctattttgc
 ccaaaTTAATTAAGGCGCGCCgctctcctgctagagtagcagtagaaaggactaccgacgaaggaact
 gggctgcgggtgtgtgattatggagtagtaagacctccctttacaacctaaaggcgaggaactgccttgcctatccaca
 atgtcgtcttcaacacattgagtcgtctccctttggaatggcccttgaccgcgccacaacctggccgcgtgaaggagtc
 caltgtctgtatttcatggtctttttacaactcatatattgtcagggtttgaaggatgcgataaggacctgtttagaca-

FIGURE 37A

Figure 1 displays 12 line drawings of insects, arranged vertically. From top to bottom, they are: a beetle (likely a ground beetle), a fly (likely a housefly), a bee (likely a honeybee), a beetle (likely a ground beetle), a fly (likely a housefly), a bee (likely a honeybee), a beetle (likely a ground beetle), a fly (likely a housefly), a bee (likely a honeybee), a beetle (likely a ground beetle), a fly (likely a housefly), and a bee (likely a honeybee). Each drawing is a detailed black and white illustration showing the insect's body, legs, and wings.

FIGURE 37B

Declaration for Patent Application

Docket Number: 1522.0030004MAC/BJD

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name.

I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter that is claimed and for which a patent is sought on the invention entitled: Compositions and Methods for Non-targeted Activation of Endogenous Genes, the specification of which is attached hereto unless the following box is checked:

- ☒ was filed on March 26, 1999;
as United States Application Number or PCT International Application Number 09/276,820; and
was amended on _____ (if applicable).

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose information that is material to patentability as defined in 37 C.F.R. § 1.56.

I hereby claim foreign priority benefits under 35 U.S.C. § 119(a)-(d) or § 365(b) of any foreign application(s) for patent or inventor's certificate, or § 365(a) of any PCT international application, which designated at least one country other than the United States listed below, and have also identified below any foreign application for patent or inventor's certificate, or PCT international application having a filing date before that of the application on which priority is claimed.

Prior Foreign Application(s)

Priority Claimed

☐ Yes ☐ No

(Application No.)

(Country)

(Day/Month/Year Filed)

☐ Yes ☐ No

(Application No.)

(Country)

(Day/Month/Year Filed)

I hereby claim the benefit under 35 U.S.C. § 119(e) of any United States provisional application(s) listed below.

(Application No.)

(Filing Date)

(Application No.)

(Filing Date)

I hereby claim the benefit under 35 U.S.C. § 120 of any United States application(s), or under § 365(c) of any PCT international application designating the United States, listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States or PCT international application in the manner provided by the first paragraph of 35 U.S.C. § 112, I acknowledge the duty to disclose information that is material to patentability as defined in 37 C.F.R. § 1.56 that became available between the filing date of the prior application and the national or PCT international filing date of this application.

09/263,814

(Application No.)

March 8, 1999

(Filing Date)

Pending

(Status – patented, pending, abandoned)

09/253,022

(Application No.)

February 19, 1999

(Filing Date)

Pending

(Status – patented, pending, abandoned)

09/159,643

(Application No.)

September 24, 1998

(Filing Date)

Abandoned

(Status – patented, pending, abandoned)

08/941,223

(Application No.)

September 26, 1997

(Filing Date)

Abandoned

(Status – patented, pending, abandoned)

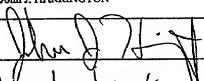
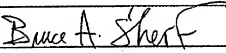
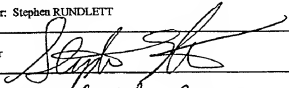
Send Correspondence to:

STERNE, KESSLER, GOLDSTEIN & FOX P.L.L.C.
1100 New York Avenue, N.W.
Suite 600
Washington, D.C. 20005-3934

Direct Telephone Calls to:

(202) 371-2600

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

Full name of sole or first inventor: John J. HARRINGTON	
Signature of sole or first inventor	 7/1/99 Date
Residence	6487 Meadowbrook Dr., Mentor OH 44060
Citizenship	USA
Post Office Address	6487 Meadowbrook Dr., Mentor OH 44060
Full name of second inventor: Bruce SHERF	
Signature of second inventor	 7/1/99 Date
Residence	7012 Avon Ct Rd, Spencer, OH, 44275
Citizenship	USA
Post Office Address	7012 Avon Ct. Rd, Spencer, OH, 44275
Full name of third inventor: Stephen RUNDLETT	
Signature of third inventor	 7/1/99 Date
Residence	703 Bell Rd. Chagrin Falls, OH, 44022
Citizenship	USA 70
Post Office Address	703 Bell Rd. Chagrin Falls, OH 44022

(Supply similar information and signature for subsequent joint inventors, if any)

Attorney's Docket No. 5817-7

PATENT

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

First Named Inventor: Harrington et al.

Group Art Unit: 1632

Application No.: 09/276,820

Examiner Name: Shukla, R.

Filed: March 26, 1999

For: COMPOSITIONS AND METHODS FOR
NON-TARGETED ACTIVATION OF
ENDOGENOUS GENES

Assistant Commissioner for Patents
Washington, DC 20231

REVOCATION OF POWER OF ATTORNEY
AND NEW POWER OF ATTORNEY BY ASSIGNEE

Sir:

Assignee hereby revokes all powers of attorney previously granted with respect to the above-identified patent application, and appoints the practitioners associated with the Customer Number provided below to prosecute this application and to transact all business in the Patent and Trademark Office connected therewith, and directs that all correspondence be addressed to that Customer Number:

Customer Number 000826

with full power of substitution and revocation to transact all business in the Patent and Trademark Office in connection therewith.

Please direct all communications to the attention of:

Anne Brown

Registration No. 36,463

Tel Raleigh Office (919) 420-2200

Fax Raleigh Office (919) 420-2260

Assignee hereby elects under 37 C.F.R. § 3.71 to prosecute this patent application and certifies that it is the assignee of the entire right, title, and interest in the patent application identified above, and in any divisionals or continuations thereof, by virtue of:

An assignment from the inventors of the patent application identified above.

The assignment was recorded in the Patent and Trademark Office at Reel 010064, Frames 0420.

In re: Harrington et al.
Appl. No.: 09/276,820
Filed: Filed: March 26, 1999
Page 2

The undersigned (whose title is supplied below) is empowered to sign this certificate on behalf of the Assignee.

ATHERSYS, INC.

By: James J. Kovach

James J. Kovach

(Print or type name of person signing)

Title: Chief Operating Officer

Date: 12/15/99

ALSTON & BIRD LLP
Post Office Drawer 34009
Charlotte, NC 28234-4009
Tel Raleigh Office (919) 420-2200
Fax Raleigh Office (919) 420-2260

CERTIFICATION OF FACSIMILE TRANSMISSION	CERTIFICATE OF MAILING
I hereby certify that this paper is being facsimile transmitted to the Patent and Trademark Office at Fax No. _____ on the date shown below. (Type or print name of person signing certification.) Signature _____ Date _____	I hereby certify that this correspondence is being deposited with the United States Postal Service with sufficient postage as first class mail in an envelope addressed to: Assistant Commissioner For Patents, Washington, DC 20231, on <u>December 30, 1999</u> . <u>Noah C. Martinez</u> NOAH C. MARTINEZ

RTA01/2071424v1

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re: Harrington, *et al.* Group Art Unit: Not Yet Assigned
Appl. No.: Not Yet Assigned Examiner: Not Yet Assigned
Filed: Filed Concurrently Herewith
For: COMPOSITIONS AND METHODS FOR NON-TARGETED ACTIVATION
OF ENDOGENOUS GENES

January 18, 2000

**REQUEST FOR TRANSFER OF COMPUTER READABLE FORM OF SEQUENCE
LISTING UNDER 37 CFR §1.821(e) AND MPEP 2422.05**

Box Patent Application
Assistant Commissioner for Patents
Washington, DC 20231

Sir:

Applicants hereby request transfer of previously filed sequence information into the above-mentioned application, concurrently filed herewith.

I hereby state that the paper copy of the sequence listing, attached hereto, is identical to the computer-readable copy of the sequence listing filed in U.S. Application Serial No. 09/276,820, filed on March 26, 1999. In accordance with 37 CFR §1.821(e) and MPEP 2422.05, please use the computer-readable form filed in that application as the computer-readable form for the above-mentioned application. It is understood that the Patent and Trademark Office will make the necessary change in application number and filing date for the present application.

Respectfully submitted,



Anne Brown
Attorney for Applicant
Registration No. 36,463

ALSTON & BIRD LLP

Post Office Drawer 34009

Charlotte, NC 28234

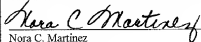
Tel Raleigh Office (919) 420-2200

Fax Raleigh Office (919) 420-2260

"Express Mail" Mailing Label Number EL247263380US

Date of Deposit: January 18, 2000

I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on the date indicated above and is addressed to Box Patent Application, Assistant Commissioner of Patents, Washington, DC 20231.



Nora C. Martinez

39

<210> 2

<211> 40

<212> DNA

<213> Homo sapiens

<400> 2

aaacttaaga tcgattaatc attcttctca tataactcaa

40

<210> 3

<211> 28

<212> DNA

<213> Homo sapiens

<400> 3

atccaccatg gctacagggtg agtactcg

28

<210> 4

<211> 36

<212> DNA

<213> Homo sapiens

<400> 4

gatccgagta ctcacctgta gccatggtgg atttaa

36

<210> 5

<211> 33

<212> DNA

<213> Homo sapiens

<400> 5

ggcgagatct agcgctatat gcgttgatgc aat

33

<210> 6

<211> 51

<212> DNA

<213> Homo sapiens

ggccagatct gctaccttaa gagagccgaa acaagcgctc atgagcccga a

51

<210> 7

<211> 6084

<212> DNA

<213> Homo sapiens

<400> 7

agatcttcaa	tattggccat	tagccatatt	attcattgggt	tatatagcat	aatcaatat	60
tggcttattg	ccattgcata	cgttgtatct	atatcataat	atgtacattt	atattggctc	120
atgtccaata	tgaccgccat	gttggcattg	attattgact	agttattaat	agtaatacat	180
tacggggcta	ttagttcata	gcccatatat	ggagttccgc	gttacataac	ttacggtaaa	240
tggcccgctc	ggctgaccgc	ccaacgaccc	cgcgccattg	acgtcaataa	tgacgtatga	300
tcccatagta	acgccaatag	ggactttcca	ttagcgtcaa	tgggtggagt	atttaccgta	360
aactgcccac	ttaggcagta	tcaaatgtga	tcatattgca	atgctcgccc	ctattgacgt	420
caatgcagat	aaatggcccc	cctggcatta	tgcccagta	atgaccttac	gggactttcc	480
tacttggcag	tacatctcat	tattagtcat	cgctattacc	atggtgatgc	ggttttggca	540
gtacaccaat	gggcgtggat	agcggtttga	ctcacgggga	ttcccaagtc	tccaccccat	600
tgacgtcaat	gggagtttgt	tttggcacca	aaatcaacgg	gactttccaa	aatgtcgtaa	660
caactgcgat	cgcccccccc	gttgacgcaa	atggggcggt	ggcggtgacg	gtggggagct	720
tatataagca	gagctcgttt	agtgaaccgt	cagatcacta	gaagctttat	tgcggttagtt	780
tatcacagtt	aaattgctaa	cgcagtcagt	gtctctgaca	caacagttct	gaacttaagc	840
tgcagtgaat	cttctaatta	actccaccag	tctcacttca	gttccttttg	cctccaccag	900
tctcacttca	gttccttttg	catgaagagc	tcagaatcaa	aagaggaaac	caacccctaa	960
gatgagcttt	ccatgtaaat	ttgtagccag	cttccttctg	attttcaatg	tttcttccaa	1020
agggtgcagtc	tccaaagaga	ttacgaatgc	cttgaaaac	tgggggtgct	tgggtcagga	1080
catcaacttg	gacatttcta	gttttcaaat	gattgatgat	attgacgata	tgaatggagc	1140
aaaaacttca	tacgaagaaa	agatttgcata	attcagaaaa	gagaaaagaga	ctttcaagga	1200
aaaagataca	gataagctat	ttaaaaattg	agctctgaaa	attaagcatc	tgaagaccga	1260
tgatcaggat	atctacaagg	tatcaatata	tgatacaaaa	ggaaaaaatg	tgttggaata	1320
aatatttgat	ttgaagattc	aagagagggg	ctcaaaacca	aagatctcct	ggacttgat	1380
caacacaacc	ctgacctgtg	aggtaatgaa	tggaaactgac	ccgaatttaa	acctgtcata	1440
agatgggaaa	catctaaaaa	tttctcagag	ggctatcaca	cacaagtgga	ccaccagctc	1500
gagtgcaaaa	tccaagtgtc	cagcagcgaa	caagtcagc	aaggaatcca	gtgtcgagcc	1560
gtgcagctgt	ccagagaaa	ggatccaggt	tgttagggcc	cgatccttct	agatcgagc	1620
tctcttaagg	tagcaaggtt	acaagacagg	ttaagggaga	ccaatagaaa	ctgggcttgt	1680

cgagacagag aagactcttg cgtttctgat aggcacctat tggctcttacg cggccgcgaa 1740
 ttccaagctt gagtattcta tcgtgtcacc taaataactt ggcgtaatca tggatcatc 1800
 tgtttctctg gtgaaattgt tatccgctca caattccaca caacatacga gccggaagca 1860
 taaagtgtaa agcctggggg gctaatagag tgagctaact cacattaatt gcgttgccgcg 1920
 atgcttccat ttgttgaggg ttaatgcttc gagaagacat gataagatac attgatgagt 1980
 ttggacaaac cacaacaaga atgcagttaa aaaaatgctt tatttgttaa atttgtgatg 2040
 ctattgcttt atttgaacc attataagct gcaataaaca agttaacaac aacaattgca 2100
 ttcattttat gtttcagggt caggggggaga tgtggggagt tttttaagc aagtaaaacc 2160
 tctacaaatg tggtaaaatc cgataaggat cgattccgga gcctgaatgg cgaatggacg 2220
 cgcctctgag cggcgctta agcgcggcgg gtgtgggtgt tacgcgcacg tgaccgctac 2280
 acttgccagc gccctagcgc ccgctccttt cgtttcttc ccttctcttc tcgccacggt 2340
 cgcggccttt ccccgtaag ctctaaatcg ggggctccct ttaggggttc gatttagtgc 2400
 tttacggcac ctcgacccca aaaaacttga ttaggggtgat ggttcacgta gtgggcatc 2460
 gccctgatag acggttttcc gccctttgac gttggagtcc acgttcttta atagtggact 2520
 cttgttccaa actggaacaa cactcaacc tatctcggtc tattcttttg atttataagg 2580
 gattttgccg atttcggcct attggttaaa aaatgagctg atttaacaaa aatttaacgc 2640
 gaattttaac aaaattataa cgtttacaat ttgcctgtg tacctcttga ggcggaaaga 2700
 accagctgtg gaattgtgtg cagttagggt gtggaaagtc ccagggtcc ccagcaggca 2760
 gaagtagtga aagcatgcat ctcaattagt cagcaaccag gtgtggaaag tccccaggct 2820
 cccagcagg cagaagtatg caaagcatgc atctcaatta gtcagcaacc atagtccgcg 2880
 cctaactcc gcccatccgc cccctaactc cgcgcagttc cgcctctct cgcgcccatg 2940
 gctgactaat tttttttatt tatgcagagg ccgaggccgc ctgcgctctt gagctattcc 3000
 agaagtagtg agggagcttt tttggaggcc taggcttttg caaaaagctt gattcttctg 3060
 acacaacagt ctcgaaacta aggctagagc caccatgatt gaacaagatg gattgcacgc 3120
 aggtttccgc gccgttgtgg tggagaggct attcgctat gactgggcac aacagacaat 3180
 cggctgctct gatccgcgcg tgttccggct gtcagcgcag gggcgcccggt ttctttttgt 3240
 caagaccgac ctgtccgggt ccttgaatga actgcaggac gaggcagcgc ggctatctg 3300
 gctggccacg accggcgctt ctgtgcgagc tgtgtctgac gttgtcactg aagcgggaag 3360
 ggactggctg ctattgggag aagtgcggcg gcaggatctc ctgtcatctc acctgtctcc 3420
 tgccgagaaa gtatccatca tggctgatgc agtgccggcg ctgcatacgc ttgatccggc 3480
 tacctgcccc ttgcaccacc aagcgaaca tcgcatcgag cgagcagcta ctggatgga 3540
 agccggctct gtcgatcagg atgatctgga cgaagagcat caggggctcg cgcagccga 3600
 actgttctcc aggtctcaagg cgcgcagctc cgacggcgag gatctcgtg tgacctatg 3660
 cgatgcctgc ttgcgaata tcatgggtgga aaatggccgc ttttctggat tcatcgactg 3720
 tggccggctg ggtgtggcg accgctatca ggacatagcg ttggctaccc gtgatattgc 3780
 tgaagagctt ggccgccaat gggctgaccg ctctctctg ctttacggta tcgccgctcc 3840
 cgattcgcag cgcctcgctt tctatcgctt tcttgacgag ttcttctgag cgggactctg 3900

gggttcgaaa tgaccgacca agcgacgccc aacctgccat cacgatggcc gcaataaaat 3960
 atctttatatt tcattacatc tgtgtgttgg ttttttgtgt gaagatcccg gtatggtgca 4020
 ctctcagtac aatctgtctt gatgccgcat agttaagcca gccccgacac ccgccaacac 4080
 ccgctgacgc gccctgacgg gcttgtctgc tcccgccatc cgcttacaga caagctgtga 4140
 ccgtctccgg gagctgcatg tgcagaggt tttcaccgct atcaccgaaa ccgcgagagc 4200
 gaaagggcct cgtgatacgc ctatttttat aggttaatat catgataata atggtttctt 4260
 agacgtcagg tggcactttt cggggaaatg tgcgcggaac ccctatttgt ttatttttct 4320
 aaatacatc aaatatgtat ccgctcatga gacaataacc ctgataaatg ctccaataat 4380
 attgaaaaag gaagagtatg agtattcaac atttccgtgt ccgcccttatt cctttttttg 4440
 cggcattttg ccttctcgtt tttgctcacc cagaacgcct ggtgaaagta aaagatgctg 4500
 aagatcagtt ggggtgcaga gtgggttaca tcgaactgga tctcaacagc ggtaagatcc 4560
 ttgagagttt tcgccccgaa gaacgttttc caatgatgag cacttttaaa gttctgctat 4620
 gtggcgcggt attatccgct attgacgccg ggcaagagca actcggtcgc cgatacact 4680
 attctcagaa tgacttggtt gactactcac cagtcacaga aaagcatctt acggatggca 4740
 tgacagtaag agaattatgc agtgctgcca taaccatgag tgataacact ggggccaact 4800
 tacttctgac aacgatcgga ggaccgaagg agctaaccgc ttttttgac aacatggggg 4860
 atcatgtaac tcgcttgatg cgttgggaac cggagctgaa tgaagccata ccaaacgacg 4920
 agcgtgacac cagctgacct gtacgaatgg caacaacgct gcgcaacta ttaactggcg 4980
 aactacttac tctagcttcc cggcaacaat taatagactg gatggaggcg gataaagttg 5040
 caggaccact tctgcgctgc gcccttcgga ctggctggtt tattgctgat aaactctggg 5100
 ccggtgagcg tgggtctcgc ggtatcattg cagcactggg gccagatggt aagccctccc 5160
 gtatcgtagt tatctacacg acgggggagtc aggcaactat ggatgaacga aatagacaga 5220
 tcgctgagat aggtgcctca ctgattaagc attggtaact gtcagacca gtttactcat 5280
 atatacttta gattgattta aaacttcatt ttaatttaa aaggatctag gtgaagatcc 5340
 tttttgataa tctcatgacc aaaatccctt aacgtgagtt ttctgtccac tgagcgtcag 5400
 acccgtgaga aaagatcaaa ggatcttctt gagatccttt ttttctgcgc gtaactctgt 5460
 gcttgcaaac aaaaaaacca ccgctaccag ccggtggttg tttgccgga caagagctac 5520
 caactctttt tccgaaggta actggcttca gcagagcgca gataccaaat actgtccttc 5580
 tagtgtagcc gtagttaggc caccacttca agaactctgt agcaccgcct acatacctcg 5640
 ctctgctaatt cotgttacca gtggctgctg ccagtgggca taagtctgtt cttaccgggt 5700
 tggactcaag acgatagtta ccggaataagg ccgagcggtc gggctgaacg gggggttcgt 5760
 gcacacgccc cagcttgagg cgaacgacct acaccgaact gagataccta cagcgtgagc 5820
 tatgagaaaag cgccacgctt ccggaaggga gaaaggcgga caggtatccg gtaagcgcca 5880
 gggtcggaac aggagagcgc acgagggagc ttccaggggg aaacgcctgg tatctttata 5940
 gtctgtcggg gtttcgccac ctctgacttg agcgtcgatt tttgtgatg tctgcagggg 6000
 ggcggagcct atggaaaaac gccagcaacg cggccttttt acggttctcg gccttttctt 6060
 ggccttttgc tcacatggct cgac 6084

<211> 6085

<212> DNA

<213> Homo sapiens

<400> 8

agatcttcaa	tattggccat	tagccatatt	attcattggt	tatatagcat	aatacaatat	60
tggtctattg	ccattgcata	cgttgtatct	atatacataat	atgtacattt	atattggctc	120
atgtccaata	tgaccgccat	gttggcattg	attcattgact	agttatataat	agtaatacat	180
tacgggggtca	ttagttcata	gcccatatat	ggagttccgc	gttacataac	ttacggtaaa	240
tggcccgctc	gggtgaccgc	ccaacgacc	cgcgccattg	acgtcaataa	tgacgtatgt	300
tcccatagta	acgccaatag	ggactttcca	ttagcgtcaa	tggggtggat	atttcaggta	360
aactgccacc	tttgcagtac	atacgaagt	tcatattgca	agtcgcgccc	ctattgacgt	420
caactgacgt	aaatggcccg	cctggcatta	tgcccagtac	atgaccttac	gggactttcc	480
tatttcggag	tacatctacg	tatttagtcat	cgctattacc	atgggtgatgc	ggttttggca	540
gtacaccaat	gggcgtggat	agcggtttga	ctcacgggga	tttccaagtc	tccaccccat	600
tgacgtcaat	gggagtttgt	tttggcacca	aaatcaacgg	gactttccaa	aatgtcgtaa	660
caactgcgat	cgcccgcccc	gttgacgcaa	atgggcggta	ggcgtgtacg	gtggggaggtc	720
tatatagaac	gagctcgttt	agtgaaccgt	cagatcacta	gaagctttat	tgcggttagt	780
tatcacagtt	aaattcgtaa	cgcagctcgt	gcttctgaca	caacagcttc	gacttaagc	840
tgcagtgact	ctcttaatta	actccaccag	tctcagttca	gttccttttg	cctccaccag	900
tctcaactca	gttctctttg	catgaagagc	tcagaatcaa	aagaggaaac	caaccctcaa	960
gatgagcttt	ccatgtaaat	ttgtagccag	cttctctctg	attttcaatg	tttcttccaa	1020
aggtgcagtc	tccaaagaga	ttacgaatgc	cttggaaac	tggggtgcct	tgggtcagga	1080
catcaacttg	gacattcccta	gtttttcaat	gagtgatgat	attgacgata	taaaatggga	1140
aaaaaacttc	gacaaagaaa	agattgcaca	attcagaaaa	gagaaagaga	ctttcaagga	1200
aaaagataca	tataagctat	ttaaaatgga	aactctgaaa	attaagcatc	tgaaagccga	1260
tgatcaggat	atctcaaaag	tatcaatata	tgatacaaaa	ggaaaaaatg	tgttggaaaa	1320
aatatttgat	ttgaagattc	aagagagggg	gtcaaaacca	aagatctcct	ggacttgtat	1380
caacacaaac	ctgacctgtg	aggtaatgaa	tggaaactgac	cccgaaattaa	acctgtatca	1440
agatggggaa	catctaaaac	tttctcagag	ggtcatcaca	cacaagtgga	ccaccagctc	1500
gagtgcaaaa	ttcaagtcca	cagcaggggaa	caaagtcagc	aaggaatcca	gtgtcgagcc	1560
tgtcagctgt	ccagagaaaag	ggatccccag	tgagtagggc	ccgatctctc	tacgtctcag	1620
ctctcttaag	gtagcaaggt	tacaagagca	gtttaaggag	accaatgaaa	actgggcttg	1680
tcgagacaga	gaagactctt	gcgtttctga	tgggacacta	tttgttctac	gcggccgcga	1740
attccaagct	tgagtattct	atcgtgtcac	ctaaataact	tggcgtaatc	atggctcatat	1800

ctgtttcctg tgtgaaattg ttatccgctc acaattccac acaacatacg agccggaagc 1860
ataaagtgtg aagcctgggg tgcctaataa gtgagctaac tcacattaat tgcgttgcgc 1920
gatgcttcca ttttgtgagg gttaatgctt cgagaagaca tgataagata cattgatgag 1980
tttgacaaa ccacaacaag aatgcagtga aaaaaatgct ttatttgtga aatttgtgat 2040
gctattgctt tatttgtaac cattataaagc tgcaataaac aagttaacaa caacaattgc 2100
attcatttta tgtttcagggt tcagggggag atgtggggagg ttttttaaag caagtaaaac 2160
ctctacaaat gtgtgtaaat ccgataagga tcgattccgg agcctgaatg gcgaatggac 2220
gcgcctgtga gcggcgcat aagcgcgagg ggtgtggtgg ttaogcgac gtgaccgcta 2280
cacttggcag cgcctagcgc ccgctcctt tcgctttctt ccttctctt ctccgacagt 2340
tcggcgctt tcccgcgcaa gctctaaatc gggggctccc tttagggttc cgatttagtg 2400
ctttacggca cctcgacccc aaaaaacttg attaggggtg tggttcacgt agtgggcat 2460
cgccctgata gacggttttt cgccttttga cgttggagtc cagcttcttt aatagtggac 2520
tcttgttcca aactggaaca aactcaacc ctatctcggg ctattctttt gatttataag 2580
ggattttgcc gatttcggcc tatttggttaa aaaatgagct gatttaacaa aaatttaacg 2640
cgaattttaa caaaatatta acgcttacaa tttcgctgt gtaccttctg aggcggaaag 2700
aaccagctgt ggaatgtgtg tcagttagggt tgtggaaagt cccagggctc cccagcaggc 2760
agaagtatgc aaagcatgca tctcaattag tcagcaacca ggtgtggaaa gtccccaggc 2820
tccccagcag gcagaagtat gcaaagcatg catctcaatt agtcagcaac catagtcccg 2880
cccctaactc cgcctatccc gccctaact ccgcccagtt ccgcccattc tccgcccatt 2940
ggctgactaa tttttttat ttatgcagag gccgaggccg cctcgccctc tgagctattc 3000
cagaagtagt gaggaggctt ttttggaggc ctaggctttt gcaaaaagct tgattcttct 3060
gacacaacag tctcgaactt aaggctagag ccacatgat tgaacaagat ggattgcacg 3120
caggttctcc ggcgccttgg gtggagaggc tattcggtca tgactgggca caacagacaa 3180
tcggctgctc tgatgccgc gtgttcgggc tgtcagcgca gggggcccg gtctttttg 3240
tcaagaccga cctgtccggg gccctgaagt aactgcagga cggggcagcg cggctatcgt 3300
ggctggccac gacgggcgtt ccttgcgcag ctgtgctcga cgttgctact gaagcgggaa 3360
gggactggct gctattgggc gaagtgccg ggcaggatct cctgtcatct cacttgcctc 3420
ctgccgagaa agtatccatc atggctgatg caatgcggcg gctgcatacg cttgatccgg 3480
ctacctgcc attcgaccac caagcgaaac atcgcatoga gcgagcacgt actcggtatg 3540
aagccggtct tgtcgatcag gatgatctgg acgaagagca tcaggggctc gcgcccaggc 3600
aactgttcgc cagggtcaaag gcgcgcatgc ccgacggcga ggatctcgtc gtgacctatg 3660
gcgatgcctg ctgtccgaat atcatggtgg aaaaaggcgg cttttctgga ttcacgact 3720
gtggccggct ggggtgtggc gaccgctatc aggacatagc gttggctacc cgtgatattg 3780
ctgaagagct tggcgcgcaa tgggctgacc gcttctcgt gctttacggg atcgccgctc 3840
ccgattcgca gcgcacgccc ttctatcgcc ttcttgacga gttcttctga gcgggactct 3900
ggggttcgaa atgaccgacc aagcgacgcc caacctgcca tcacgatggc cgcaataaaa 3960
tatctttatt ttcattacat ctgtgtgttg gtttttgtg tgaagatccg cgtatggtgc 4020

actctcagta caatctgctc tgatgccgca tagttaagcc agccccgaca cccgcccaaca 4080
 ccgcctgacg cgcctctgacg ggcttgctcg ctcccggcat cgccttacag acaagctgtg 4140
 accgtctccg ggagctgcac gtgtcagagg ttttcaccgt catcaccgaa acgcgcgaga 4200
 cgaaagggcc tcgtgatacg cctattttta taggttaatg tcattgataat aatggtttct 4260
 tagacgtcag gtggcacttt tcggggaaat gtgcgcggaa cccctatttg tttatttttc 4320
 taaatacatt caaatatgta tcgcctcatg agacaataac cctgataaat gcttcaataa 4380
 tattgaaaaa ggaagagtat gagtattcaa catttcogtg tcgcccttat tccctttttt 4440
 gcggcatttt gccttctctg ttttgctcac ccagaaacgc tggtgaaagt aaaagatgct 4500
 gaagatcagt tgggtgcacg agtgggttac atcgaaactg atctcaacag cggtaagatc 4560
 cttgagagt ttcgccccga agaacgtttt ccaatgatga gcacttttaa agttctgcta 4620
 tgtggcgcgg tattatcccg tattgacgcc gggcaagagc aactcggtcg ccgcatacac 4680
 tattctcaga atgacttggt tgagtactca ccagtcacag aaaagcatct tacggatggc 4740
 atgacagtaa gagaattatg cagtgtctgcc ataaccatga gtgataacac tgcggccaac 4800
 ttacttctga caacgatcgg aggaccgaag gagctaacgg cttttttgca caacatgggg 4860
 gatcatgtaa ctgccttga tcgttgggaa ccggagctga atgaagccat accaaaacgc 4920
 gagcgtgaca ccaecatgac tgtagcaatg gcaacaacgt tgcgcaaat attaaactgc 4980
 gaactactta ctctagcttc ccggcaacaa ttaatagact ggtatggagg ggataaaagt 5040
 gcaggaccac ttctgcgctc ggcccttcog gctggctggt ttattgctga taaatctgga 5100
 gccggtgagc gtgggtctcg cgggtatcatt gcagcactgg gccagatgg taagccctcc 5160
 cgtatctgag ttatctacac gacggggagt caggcaacta tggatgaacg aaatagacag 5220
 atcgctgaga taggtgcctc actgattaa ctttggtaac tgtcagacca agtttactca 5280
 tatatacttt agattgattt aaaaactcat ttttaattta aaaggatcta ggtgaagatc 5340
 ctttttgata atctcatgac caaaatccct taacgtgagt tttcgttcca ctgagcgtca 5400
 gaccocgtag aaaagatcaa aggatctctt tgagatcctt ttttctcg cgtaactcgc 5460
 tgcttgcaaa caaaaaaac accgctacca gcggtggttt gtttcccgga tcaagagcta 5520
 ccaactcttt ttccgaaggt aactggcttc agcagagcgc agataccaaa tactgtcctt 5580
 ctagtgtagc cgtagttagg ccaccacttc aagaactctg tagcaccgcc tacatacctc 5640
 gctctgctaa tctgttacc agtggctgct gccagtgcg ataagtcgtg tcttaccggg 5700
 ttggactcaa gacgatagtt accggataag gcgcagcgg cgggctgaac ggggggttcg 5760
 tgcacacagc ccagcttgga ggaacgacc tacaccgaac tgagatacct acagcgtgag 5820
 ctatgagaaa gcgccacgct tcccgaaggg agaaaaggcg acaggtatcc ggtaaagcgc 5880
 agggctcgaa caggagagcg cacgagggag cttccagggg gaaacgcctg gtatctttat 5940
 agtcctgtcg ggtttcgcca cctctgactt gagcgtcgat ttttgtgatg ctgcgcaggg 6000
 gggcgaggcc tatggaaaaa ccgcagcaac gcggcctttt tacggttctt ggccttttgc 6060
 tggccttttg ctcaactggc tcgac 6085

<211> 6086

<212> DNA

<213> Homo sapiens

<400> 9

```

agatcttcaa tattggccat tagccatatt attcattggt tatatagcat aaatcaatat 60
tggctattgg ccattgcata cgttgtatct atatcataat atgtacattt atattggctc 120
atgtccaata tgacggccat gttggcattg attattgact agttattaat agtaatacat 180
tacgggggtca ttagttcata gcccatatat ggagttccgc gttacataac ttacggtaaa 240
tggcccgctt ggctgaccgc ccaacgaccc ccgcccattg acgtcaataa tgacgtatgt 300
tcccatagta acgccaatag ggactttcca ttgacgtcaa tgggtggagt atttacggta 360
aactgccccac ttggcagtag atcaagtgtg tcatatgcc agtccgcccc ctattgacgt 420
caatgacggt aaatggcccg cctggcatta tgcccagtag atgaccttac gggactttcc 480
tacttggcag tacatctacg tattagtcac cgtattacc atgggtgtagt ggttttggca 540
gtacaccaat gggcgtaggt agcggtttga ctcacgggga ttccaagtc tccaccccat 600
tgacgtcaat gggagtttgt tttggcacca aaatcaacgg gactttccaa aatgtcgtaa 660
caactgcgat cgcgcgcccc gttgacgcaa atgggcggta ggctgtacg gtgggaggtc 720
tatataagca gacgtcggtt agtgaaccgt cagatcacta gaagctttat tgcggtaggt 780
tatcacagtt aaattgctaa cgcagtcagt gcttctgaca caacagtcct gaacttaagc 840
tgacgtgact ctcttaatta actccaccag tctcacttea gttccttttg cctccaccag 900
tctcacttea gttccttttg catgaagagc tcagaatcaa aagaggaaac caaccctaa 960
gatgagcttt ccatgtaaat ttgtagccag ctctcctctg attttcaatg tttcttccaa 1020
aggtgacgtc tccaaagaga ttacgaatgc cttggaiaac tggggtgcct tgggtcagga 1080
catcaacttg gacattccta gttttcaaat gaggtagtat attgacgata taaaatggga 1140
aaaaacttca gacaagaaaa agattgcaca attcagaaaa gagaaagaga ctttcaagga 1200
aaaagataca tataagctat ttaaaatagg aactctgaaa attaagcatc tgaagaccga 1260
tgatcaggat atctacaagg tatcaatata tgatacaaaa ggaaaaaatg tgttggaaaa 1320
aatatttgat ttgaagattc aagagagggg ctcaaaaacca aagatctcct ggaactgtat 1380
caacacaacc ctgacctgtg aggtaaatga ttgaactgac ccgaattaa acctgtatca 1440
agatgggaaa catctaaaac tttctcagag ggtcatcaca cacaagtgga ccaccagcct 1500
gagtgcaaaa ttcaagtga cagcagggaa caaagtcagc aaggaaatcca gtgtcgagcc 1560
tgtcagctgt ccagagaaag ggtccacag gtgagtaggg ccgatcctt ctgagatcga 1620
gctctcttaa ggtagcaagg ttacaagaca ggtttaagga gaccaataga aactgggctt 1680
gtcgagacag agaagactct tgcttttctg ataggcacct attggtctta cgcggccgcg 1740
aattccaagc ttgagtattc tctcgtgtca cctaaataac ttggcgtaat catggtcata 1800
tctgtttcct ggtgaaatt gtatccgctt cacaattcca cacaacatag gaggcggaag 1860
cataaagtgt aaagcctggg gtgcctaagt agtgagctaa ctcacattaa ttgcgttgcg 1920

```

cgatgcttcc attttgtgag ggttaatgct tgcagaagac atgataagat acattgatga 1980
 gtttggacaa accacacaacaa gaatgcagtg aaaaaaatgc ttattttgtg aaattttgtga 2040
 tgctatttgc ttattttgtaa ccattataag ctgcaataaa caagttaaca acaacaattg 2100
 cattcatttt atgttttcagg ttcaggggga gatgtgggag gttttttaaa gcaagtaaaa 2160
 cctctacaaa tgtggttaaa tccgataagg atcgattccg gagcctgaat ggcgaatgga 2220
 cgcgcctgt agcggcgcat taagcgcggc ggggtgtggt gttacgcgca cgtgaccgct 2280
 acacttgcca gcgccttagc gccgcctcct ttcgctttct tcccttcctt tctgcgca 2340
 ttccggcgtc ttcccgctca agctctaaat cgggggctcc ctttagggtt ccgatttagt 2400
 gctttacggc acctcgacc caaaaaactt gattagggtg atggttcacg tagtgggcca 2460
 tcgcccctgat agacggtttt tgcgcccttg acgttggagt ccacgttctt taatagtga 2520
 ctcttgcttc aaactggtaac aacactcaac cctatctcgg tctattcttt tgatttataa 2580
 gggattttgc cgatttcggc ctattgggta aaaaatgagc tgatttaaca aaatttaac 2640
 gcgaatttta acaaaaatatt aacgcttaca atttcgcctg tgtaccttct gagggcgaaa 2700
 gaaccagctg tggaatgtgt gtcagttagg gtgtggaaa tccccaggt cccagcagg 2760
 cagaagtatg caaagcatgc atctcaatta gtcagcaacc aggtgtggaa agtccccagg 2820
 ctccccagca ggcagaagta tgcaaagcat gcactcfaat tagtcagcaa ccatagtccc 2880
 gccctaaact cgcgccatcc gcgccctaac tccgccagc tccgccctc ctccgcccca 2940
 ttggtgacta atttttttta ttatgcaga ggcgcaggcc gcctcggcct ctgagctatt 3000
 ccagaagtat tgaggaggct tttttggagg cctaggcttt tgcaaaaagc ttgattcttc 3060
 tgacacaaca gtttcgaact taaggctaga gccaccatga ttgaacaaga tggattgcac 3120
 gcaggtttct cggccgcttg ggtggagagg ctattcggct atgactgggc acaacagaca 3180
 atcggctgct ctgatgcgc cgtgttccgg ctgtcagcgc aggggcgcgc ggtctctttt 3240
 gtcaagaccg acctgtccgg tgcctgaat gaactgcagg acgaggcagc gcggctatcg 3300
 ttgctggcca cgacggcgct tccttgcgca gctgtgctgc acgttgcac tgaagcggga 3360
 agggactggc tgctattggg cgaagtgcgc ggcgaggatc tcctgtcac tcacctgct 3420
 cctgcgcaga aagtatccat catggctgat gcaatgcgc gctgcatac gcttgatcg 3480
 gctacctgcc cattcgacca ccaagcgaaa catcgcatgc agcagcagc tactcggtat 3540
 gaagccggtc ttgtcgatca ggtgatctg gacgaagagc atcaggggct cgcgccagcc 3600
 gaactgttgc ccaggctcaa ggcgcgcgat cccgacggcg aggatctcgt cgtgacccat 3660
 ggcgatgctt gcttgcgcaa tatcatggtg gaaaaatggc gcttttcttg attcatcgac 3720
 tgtggccggc tgggtgtggc ggaccgctat caggacatag cgttggttac ccgtgatatt 3780
 gctgaagagc ttggcggcga atgggctgac cgcttcctcg tgccttacgg tatcgccgct 3840
 ccgatttcgc agcgcategc ctctcatcgc ctctctctcg agttctctcg agcgggactc 3900
 ttgggttcga aatgaccgac caagcgacgc ccaacctgcc atcagcatg ccgcaataaa 3960
 atatctttat ttctattaca tctgtgtgtt ggttttttgt gtgaagatcc gcgtatggtg 4020
 cactctcagt acaatctgct ctgatgcgc atagttaagc cagccccgac acccgccaac 4080
 acccgctgac gcgcctgac gggctgtgct gctccggca tccgcttaca gacaagctgt 4140

gaccgtctcc gggagctgca tgtgtcagag gttttcaccg tcatcaccga aacgcgcgag 4200
 acgaaagggc ctctgtatgc gctattttt ataggtaaat gtcagtataa taatgggttc 4260
 ttagacgtca ggtggcactt ttccgggaaa tgtgcgcgga acccctattt gtttattttt 4320
 ctaaatacat tcaaatatgt atccgctcat gagacaataa ccttgataaa tgcttcaata 4380
 atattgaaaa aggaagagta ttagtattca acatttcctg gtcgccctta ttcctttttt 4440
 tgcgcgcat tgccttctctg tttttgctca ccagaaacg ctgggtgaaag taaaagatgc 4500
 tgaagatcag ttgggtgcac gagtgggtta catcgaactg gatctcaaca gcggtgaagat 4560
 ccttgagagt ttgcgcccg aagaacgttt tccaatgatg agcactttta aagttctgct 4620
 atgtggcgcg gtattatccc gtattgacgc cgggcaagag caactcggtc gccgcataca 4680
 ctattctcag aatgacttgg ttgagtactc accagtcaca gaaaagcatt ttacggatgg 4740
 catgacagta agagaattat gcagtgtctgc cataaccatg agtgataaca ctgcggccaa 4800
 ctacttctg acaacgatcg gaggaaccga ggagctaacc gcttttttgc acaacatggg 4860
 ggatcatgta actcgccttg atcgttggga accggagctg aatgaagcca taccaaaacga 4920
 cgagcgtgac accacgatgc ctgtagcaat ggcaacaacg ttgcgcaaac tattaactgg 4980
 cgaactactt actctagctt ccgcggcaaca attaatagac tggatggagg cggataaagt 5040
 tgcaggacca ctctcgcct cggcccttcc ggctggctgg tttattgctg ataatctgg 5100
 agccggtgag cgtgggtctc gcggtatcat tgcagcactg gggccagatg gtaagccctc 5160
 ccgtatcgta gttatctaca cgacggggag tcaggcaact atggatgaac gaaatagaca 5220
 gatcgtgag atagggtgct cactgattaa gcattggtta ctgtcagacc aagtttactc 5280
 atataactt tagattgatt taaaacttca tttttaattt aaaaggatct aggtgaagat 5340
 cctttttgat aatctcatga ccaaaatccc ttaacgtgag ttttcgttcc actgagcgtc 5400
 agaccccgta gaaaagatca aaggatcttc ttgagatcct tttttcttgc gcgtaactcg 5460
 ctgcttgcaa acaaaaaaac caccgctacc agcgttggtt tgtttgcggg atcaagagct 5520
 accaactctt ttccgaagg taactggctt cagcagagcg cagataccaa atactgtcct 5580
 tctagtgtag ccgtagttag gccaccactt caagaactct gtacacgcgc ctacatacct 5640
 cgctctgcta atcctgttac cagtggctgc tgcagtgagg gataagtcgt gtcttaccgg 5700
 gttggactca agacgatagt tacgggataa ggcgcagcgg tcgggctgaa cgggggggtc 5760
 gtgcacacag ccagcgttgg agcgaacgac ctacaccgaa ctgagatacc tacagcgtga 5820
 gctatgagaa agcgcacgcg ttcccgaagg gagaaggcgg gacaggtatc cggtaagcgg 5880
 cagggtcgga acaggagagc gcacgagggg gcttccaggg ggaacgcct ggtatcttta 5940
 tagtctgtc gggtttccgc acctctgact tgagcgtcga tttttgtgat gctcgtcagg 6000
 ggggcgcgagc ctatggaaaa acgccagcaa cgcgcgcctt ttacgggttc tggccttttg 6060
 ctggcctttt gctcacatgg ctcgac 6086

<210> 10

<211> 38

<212> DNA

<213> Artificial sequence

<220>

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 10

tttttttttt ttgcgcagcg gcgcacatcnn nntttatt 38

<210> 11

<211> 25

<212> DNA

<213> Artificial sequence

<220>

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 11

cagatcacta gaagctttat tgcgg 25

<210> 12

<211> 20

<212> DNA

<213> Artificial sequence

<220>

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 12

tttcgctcag cggccgcac 20

<210> 13

<211> 45

<212> DNA

<213> Artificial sequence

<220>

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 13

actcataggc catagaggcc tatcacagtt aaattgctaa cgcag

45

<210> 14

<211> 43

<212> DNA

<213> Artificial sequence

<221> OTHER

<222> 1

<223> 5' cytosine at position #1 is biotinylated

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 14

ctcgtttagt gcggccgctc agatcactga attctgacga cct

43

<210> 15

<211> 41

<212> DNA

<213> Artificial sequence

<221> OTHER

<222> 1

<223> 5' cytosine at position #1 is biotinylated

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 15

ctcgtttagt gcgcgccag atcactgaat tctgacgacc t

41

<210> 16

<211> 22

<212> DNA

<213> Artificial sequence

<221> OTHER

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 16

gacctactga ttaacggcca ta

22

<210> 17

<211> 20

<212> DNA

<213> Artificial sequence

<221> OTHER

<222> 1

<223> 3' thymidine at position #20 is biotinylated

<223> Description of artificial sequence: synthetic oligonucleotide

<400> 17

tcgtcagaat tcagtgatct

20